

# The Impact of Sleep Patterns on Body Mass Index: Study from NHANES Dataset

Prathamesh Joshi

## Background

Sleep quality has been recognized as a critical component of overall health. Insufficient or poor-quality sleep (e.g., less than 6 hours of sleep, frequent awakenings) has been linked to numerous health issues, including obesity and hypertension. This study focuses on understanding how sleep patterns, including sleep duration, snoring frequency, and daytime sleepiness, relate to Body Mass Index (BMI). By analyzing data from the NHANES dataset, this study aims to provide insights that could contribute to lifestyle interventions and public health strategies aimed at improving sleep cycles and managing health risks associated with poor sleep.

## Data Source

Data will be extracted from the **National Health and Nutrition Examination Survey (NHANES)**.

## Study Parameters

The study focuses on adults, using data on sleep patterns (e.g., sleep duration, snoring frequency, daytime sleepiness), BMI, and demographic variables such as age, gender, and ethnicity.

So far, the following parameters have been identified from the NHANES dataset for this analysis:

### 1. Sleep Patterns:

- **Sleep duration (SLD012):** Usual sleep time on weekdays
- **Usual sleep time on weekdays (SLQ300)**
- **Usual wake time on weekdays (SLQ310)**
- **Snoring frequency (SLQ030)**
- **Daytime sleepiness (SLQ120)**

### 2. Health Outcomes:

- **BMI (BMXBMI)**

### 3. Demographic Variables:

- **Age (RIDAGEYR)**
- **Gender (RIAGENDR)**
- **Ethnicity (RIDRETH1)**

## Study Aims

# Primary Study Aim

We want to study how sleep patterns (like how long people sleep, how often they snore, and how sleepy they feel during the day) relate to BMI in adults.

## Secondary Study Aims

- We will check if demographic factors (like age, gender, and ethnicity) affect the link between sleep behaviors and health outcomes.
- We will look at how often different sleep quality indicators appear in various demographic groups.

# Code Structure

The code used to analyze the NHANES dataset has been broken into two sections to enhance clarity and flow:

## 1. Function Definitions Section

This section contains all the necessary functions that will be used in the code execution section. It's important for you to review and understand the function definitions, as they provide the foundation for the subsequent analysis.

## 2. Code Execution Section

This is where the actual data analysis takes place. After you have gone through the function definitions, you can scroll down to this section to see how the functions are applied to the data and how the results are derived.

It's best to scroll through the function definitions and move on to the code execution section to proceed with the report.

## 1. Function Definitions Section

```
if (!require(nhanesA)) install.packages("nhanesA")
```

```
## Loading required package: nhanesA
```

```
if (!require(dplyr)) install.packages("dplyr")
```

```
## Loading required package: 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
```

```
if (!require(tidyr)) install.packages("tidyr")
```

```
## Loading required package: tidyr
```

```
if (!require(lmtest)) install.packages("lmtest")
```

```
## Loading required package: lmtest
```

```
## Loading required package: zoo
```

```
##  
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':  
##  
## as.Date, as.Date.numeric
```

```
if (!require(ggplot2)) install.packages("ggplot2")
```

```
## Loading required package: ggplot2
```

```
if (!require(car)) install.packages("car")
```

```
## Loading required package: car
```

```
## Loading required package: carData
```

```
##  
## Attaching package: 'car'
```

```
## The following object is masked from 'package:dplyr':  
##  
## recode
```

```
if (!require(ggcorrplot)) install.packages("ggcorrplot")
```

```
## Loading required package: ggcorrplot
```

```
if (!require(haven)) install.packages("haven") # For reading XPT files
```

```
## Loading required package: haven
```

```
# Load libraries
library(nhanesA)
library(dplyr)
library(tidyr)
library(ggplot2)
library(ggcorrplot)
library(haven)
library(lmtest)
library(car) # For vif()
library(gridExtra)
```

```
##
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
##
##      combine
```

```

# Function to clean and merge datasets
clean_data <- function(demo, slq, bmx) {
  # Merge the datasets by SEQN (Survey Participant ID)
  merged_data <- demo %>%
    left_join(slq, by = "SEQN") %>%
    left_join(bmx, by = "SEQN")

  # Ensure the columns are present before filtering
  essential_columns <- c("BMXBMI", "SLD012", "SLQ030", "SLQ120", "SLQ310", "SLQ300")
  missing_columns <- setdiff(essential_columns, names(merged_data))

  if(length(missing_columns) > 0) {
    stop(paste("Missing columns:", paste(missing_columns, collapse = ", ")))
  }

  # Clean the data: remove rows with missing essential values
  cleaned_data <- merged_data %>%
    filter(!is.na(BMXBMI), !is.na(SLD012), !is.na(SLQ030), !is.na(SLQ120), !is.na(SLQ310), !is.na(SLQ300)) %>%
    dplyr::select(SEQN, RIDAGEYR, RIAGENDR, RIDRETH1, SLD012, SLQ030, SLQ120, SLQ310, SLQ300, BMXBMI) %>%
    # Convert categorical variables to factors
    mutate(
      RIAGENDR = factor(RIAGENDR, levels = c(1, 2), labels = c("Male", "Female")),
      RIDRETH1 = factor(RIDRETH1, levels = c(1, 2, 3, 4, 5),
        labels = c("Non-Hispanic White", "Non-Hispanic Black", "Mexican American",
          "Other Hispanic", "Other Race")),
      SLD012 = as.numeric(SLD012), # Keep SLD012 as numeric for sleep hours

      # SLQ030 - How often do you snore?
      SLQ030 = factor(SLQ030, levels = c(0, 1, 2, 3, 7, 9), labels = c("Never", "Rarely - 1-2 nights a week",
        "Occasionally - 3-4 nights a week", "Frequently - 5 or more nights a week", "Refused", "Don't know")),

      # SLQ120 - How often feel overly sleepy during day?
      SLQ120 = factor(SLQ120, levels = c(0, 1, 2, 3, 4, 7, 9), labels = c("Never", "Rarely - 1 time a month",
        "Sometimes - 2-4 times a month", "Often - 5-15 times a month", "Almost always - 16-30 times a month", "Refused", "Don't know")),
    )
}

```

```

SLQ310 = as.character(SLQ310), # Treat SLQ310 as a character variable ('HH:MM')
SLQ300 = as.character(SLQ300), # Treat SLQ300 as a character variable ('HH:MM')

# Validate time format for SLQ310 and SLQ300
SLQ310 = ifelse(grepl("^\\d{2}:\\d{2}$", SLQ310), SLQ310, NA),
SLQ300 = ifelse(grepl("^\\d{2}:\\d{2}$", SLQ300), SLQ300, NA),

# Assign sleep categories based on SLQ300 (sleep time) and SLQ310 (wake time)
SLQ300_category = factor(
  case_when(
    as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) >= 19
* 60 & as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) < 21 * 6
0 ~ "SLEPT EARLY",
    as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) >= 21
* 60 & as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) < 23 * 6
0 ~ "SLEPT ON TIME",
    as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) >= 23
* 60 | as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) < 1 * 60
~ "SLEPT LATE", # Between 23:00 and 01:00
    as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) >= 1
* 60 & as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) < 3 * 60
~ "SLEPT VERY LATE",
    as.numeric(substr(SLQ300, 1, 2)) * 60 + as.numeric(substr(SLQ300, 4, 5)) >= 3
* 60 ~ "SLEPT AT OTHER TIME",
    TRUE ~ "NA"
  ),
  levels = c("SLEPT EARLY", "SLEPT ON TIME", "SLEPT LATE", "SLEPT VERY LATE", "SLE
PT AT OTHER TIME")
),

SLQ310_category = factor(
  case_when(
    as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) >= 4
* 60 & as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) < 6 * 60
~ "WOKE UP EARLY",
    as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) >= 6
* 60 & as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) < 8 * 60
~ "WOKE UP ON TIME",
    as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) >= 8
* 60 & as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) < 10 * 6
0 ~ "WOKE UP LATE",
    as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) >= 10
* 60 & as.numeric(substr(SLQ310, 1, 2)) * 60 + as.numeric(substr(SLQ310, 4, 5)) < 12 * 6
0 ~ "WOKE UP VERY LATE",
    TRUE ~ "WOKE UP AT OTHER TIMES"
  ),
  levels = c("WOKE UP EARLY", "WOKE UP ON TIME", "WOKE UP LATE", "WOKE UP VERY LAT
E", "WOKE UP AT OTHER TIMES")
)
) %>%

# Optional: Additional cleaning - Remove rows with missing 'BMXBMI' values

```

```

    filter(!is.na(BMXBMI))

  return(cleaned_data)
}

# Function to calculate descriptive statistics for cleaned data
descriptive_stats <- function(cleaned_data) {
  # Descriptive statistics for numerical variables
  numeric_vars <- cleaned_data %>%
    select(where(is.numeric))

  numeric_stats <- numeric_vars %>%
    summarise(across(everything(), list(
      mean = ~mean(. , na.rm = TRUE),
      sd = ~sd(. , na.rm = TRUE),
      min = ~min(. , na.rm = TRUE),
      max = ~max(. , na.rm = TRUE)
    )))

  # Descriptive statistics for categorical variables
  categorical_vars <- cleaned_data %>%
    select(where(is.factor) | where(is.character))

  categorical_stats <- categorical_vars %>%
    summarise(across(everything(), ~{
      counts <- table(.)
      paste(names(counts), counts, sep = ": ", collapse = ", ")
    }))

  # Combine both numeric and categorical stats
  stats <- list(
    numeric = numeric_stats,
    categorical = categorical_stats
  )

  return(stats)
}

# Correlation Analysis Function
correlation_analysis <- function(cleaned_data) {
  # Select only numerical columns for correlation analysis (excluding SEQN)
  numeric_vars <- cleaned_data %>%
    select(where(is.numeric)) %>%
    select(-SEQN) # Exclude SEQN from correlation analysis

  # Compute the correlation matrix
  correlation_matrix <- cor(numeric_vars, use = "complete.obs", method = "pearson")

  # Return the correlation matrix
  return(correlation_matrix)
}

```

```

# Function for multiple linear regression
regression_analysis <- function(cleaned_data) {

  # Create the regression model
  model <- lm(BMXBMI ~ SLD012 + SLQ030 + SLQ120 + RIDAGEYR + RIAGENDR, data = cleaned_data)

  # Summary of the model
  model_summary <- summary(model)
  # Step 2: Visualize Residuals
  # 1) Residuals vs Fitted plot
  residuals_plot <- ggplot(data = data.frame(fitted = model$fitted.values, residuals = model$residuals),
                           aes(x = fitted, y = residuals)) +
    geom_point() +
    geom_smooth(method = "loess", se = FALSE, color = "blue") + # Loess smooth line
    labs(title = "Residuals vs Fitted", x = "Fitted Values", y = "Residuals") +
    theme_minimal()

  # 2) Q-Q Plot of residuals
  qq_plot <- ggplot(data = data.frame(residuals = model$residuals), aes(sample = residuals)) +
    stat_qq() +
    stat_qq_line() +
    labs(title = "Q-Q Plot of Residuals") +
    theme_minimal()

  # 3) Fitted Values vs sqrt(Standardized Residuals)
  standardized_residuals <- rstandard(model)
  sqrt_standardized_residuals <- sqrt(abs(standardized_residuals))

  fitted_vs_sqrt_residuals_plot <- ggplot(data = data.frame(fitted = model$fitted.values, sqrt_residuals = sqrt_standardized_residuals),
                                           aes(x = fitted, y = sqrt_residuals)) +
    geom_point() +
    geom_smooth(method = "loess", se = FALSE, color = "blue") + # Loess smooth line
    labs(title = "Fitted Values vs Sqrt(Standardized Residuals)", x = "Fitted Values", y = "Sqrt(Standardized Residuals)") +
    theme_minimal()

  # 4) Leverage vs Standardized Residuals plot
  leverage_values <- hatvalues(model)
  leverage_vs_residuals_plot <- ggplot(data = data.frame(leverage = leverage_values, standardized_residuals = standardized_residuals),
                                           aes(x = leverage, y = standardized_residuals)) +
    geom_point() +
    geom_smooth(method = "loess", se = FALSE, color = "blue") + # Loess smooth line
    labs(title = "Leverage vs Standardized Residuals", x = "Leverage", y = "Standardized Residuals") +
    theme_minimal()

  # Step 3: Check for Multicollinearity (Variance Inflation Factor)

```



```

vif_values <- vif(model)

# Display VIF values
print("Variance Inflation Factor (VIF) values:")
print(vif_values)

# Return model summary and diagnostics
return(list(
  model_summary = model_summary,
  residuals_plot = residuals_plot,
  qq_plot = qq_plot,
  fitted_vs_sqrt_residuals_plot = fitted_vs_sqrt_residuals_plot,
  leverage_vs_residuals_plot = leverage_vs_residuals_plot,
  vif_values = vif_values
))
}

# Function for advanced intereaction effect analysis
# I am truncating this for Rmd for formating purposes, since the code is a little long,
# Please refer to the
# attached .R file for complete code

```

## 2. Function Execution Section

### Step 1: Load the Datasets

I started by loading the three datasets: DEMO\_I.XPT , SLQ\_I.XPT , and BMX\_I.XPT . These datasets contain different types of information about survey participants.

```

# Step 1: Load the datasets
demo <- read_xpt("/Users/pratham/project_aston/fall_2024/CSnP/DEMO_I.XPT")
slq <- read_xpt("/Users/pratham/project_aston/fall_2024/CSnP/SLQ_I.XPT")
bmx <- read_xpt("/Users/pratham/project_aston/fall_2024/CSnP/BMX_I.XPT")

```

### Step 2: Clean and Merge the Data

Next, I cleaned and merge the data into one dataset using the clean\_data() function. This function does several important things:

## 2.1 Merging the Datasets

## 2.2 Checking for Missing Columns

## 2.3 Filtering Out Missing Values

## 2.4 Converting Variables to Factors

## 2.5 Cleaning Time Data and Assigning Categories

So basically in this function I have cleaned the data to ensure it was complete, removing any rows with missing values for essential variables and ensuring the data types were correctly formatted. This allows for accurate analysis and avoids errors during computation.

```
# Step 2: Clean and merge the data  
cleaned_data <- clean_data(demo, slq, bmx)
```

## Step 3) Descriptive Statistics

In this step, I calculated the descriptive statistics for both numerical and categorical variables in the cleaned data. This helps summarize key measures like mean, standard deviation, minimum, and maximum for numerical data, and frequency distributions for categorical data, providing a quick overview of the dataset.

Taking a look at it I can see the summary of the cleaned data reveals a diverse dataset with a range of ages (16 to 80 years) and BMI values (14.5 to 67.3). It shows a near-equal distribution of genders, with the majority identifying as either non-Hispanic White or Mexican American, and provides detailed distributions for sleep-related variables, with varying frequencies of sleep and wake times.

```
# Step 3: Descriptive statistics for BMI  
stats <- descriptive_stats(cleaned_data)  
print("Descriptive Statistics for BMI:")
```

```
## [1] "Descriptive Statistics for BMI:"
```

```
print(stats)
```

```
## $numeric
## # A tibble: 1 × 16
##   SEQN_mean SEQN_sd SEQN_min SEQN_max RIDAGEYR_mean RIDAGEYR_sd RIDAGEYR_min
##   <dbl>    <dbl>    <dbl>    <dbl>          <dbl>        <dbl>        <dbl>
## 1   88699.   2877.    83732   93702          46.3         19.3         16
## # i 9 more variables: RIDAGEYR_max <dbl>, SLD012_mean <dbl>, SLD012_sd <dbl>,
## #   SLD012_min <dbl>, SLD012_max <dbl>, BMXBMI_mean <dbl>, BMXBMI_sd <dbl>,
## #   BMXBMI_min <dbl>, BMXBMI_max <dbl>
##
## $categorical
## # A tibble: 1 × 8
##   RIAGENDR RIDRETH1 SLQ030 SLQ120 SLQ300_category SLQ310_category SLQ310 SLQ300
##   <chr>    <chr>    <chr>  <chr>  <chr>          <chr>          <chr>  <chr>
## 1 Male: 28... Non-His... Never... Never... SLEPT EARLY: 2... WOKE UP EARLY:... 00:00... 00:00...
```

```
print(summary(cleaned_data))
```

```

##          SEQN          RIDAGEYR          RIAGENDR          RIDRETH1
## Min.      :83732  Min.      :16.00  Male :2858  Non-Hispanic White:1065
## 1st Qu.:86190  1st Qu.:29.00  Female:3101  Non-Hispanic Black: 777
## Median :88706  Median :46.00                      Mexican American  :1907
## Mean    :88699  Mean    :46.26                      Other Hispanic   :1267
## 3rd Qu.:91196  3rd Qu.:62.00                      Other Race        : 943
## Max.     :93702  Max.     :80.00
##
##          SLD012                      SLQ030
## Min.      : 2.00  Never                      :1688
## 1st Qu.: 7.00  Rarely - 1-2 nights a week      :1404
## Median : 8.00  Occasionally - 3-4 nights a week :1005
## Mean    : 7.74  Frequently - 5 or more nights a week:1440
## 3rd Qu.: 8.50  Refused                      : 2
## Max.     :14.50  Don't know                      : 420
##
##                                     SLQ120          SLQ310
## Never                                     :1081  Length:5959
## Rarely - 1 time a month                  :1380  Class :character
## Sometimes - 2-4 times a month            :1933  Mode  :character
## Often - 5-15 times a month               :1071
## Almost always - 16-30 times a month: 489
## Refused                                 : 0
## Don't know                             : 5
##          SLQ300          BMXBMI          SLQ300_category
## Length:5959          Min.      :14.50  SLEPT EARLY          : 281
## Class :character    1st Qu.:24.00  SLEPT ON TIME         :2623
## Mode  :character    Median :28.00  SLEPT LATE            :2314
##                      Mean    :29.16  SLEPT VERY LATE       : 464
##                      3rd Qu.:32.80  SLEPT AT OTHER TIME: 277
##                      Max.     :67.30
##
##                                     SLQ310_category
## WOKE UP EARLY          :1569
## WOKE UP ON TIME        :2941
## WOKE UP LATE           : 838
## WOKE UP VERY LATE      : 262
## WOKE UP AT OTHER TIMES: 349
##
##

```

## Step 4) Correlation Analysis

The correlation analysis was performed to explore the relationships between sleep patterns and BMI. By calculating the correlation matrix, I can identify any potential linear associations between variables like age, sleep duration (SLD012), and BMI (BMXBMI).

The correlation matrix shows a weak positive correlation (0.10) between age (RIDAGEYR) and BMI (BMXBMI), indicating a slight increase in BMI with age. Sleep duration (SLD012) shows weak negative correlations with both age and BMI, suggesting that sleep duration has minimal impact on these variables in this dataset.

This is just how I got started because I had a lot of categorical variables related to sleep as well where in we might actually see some fruitful results

```
# Step 4: Correlation analysis to check relationships between sleep patterns and BMI
cor_matrix <- correlation_analysis(cleaned_data)
print("Correlation Matrix:")
```

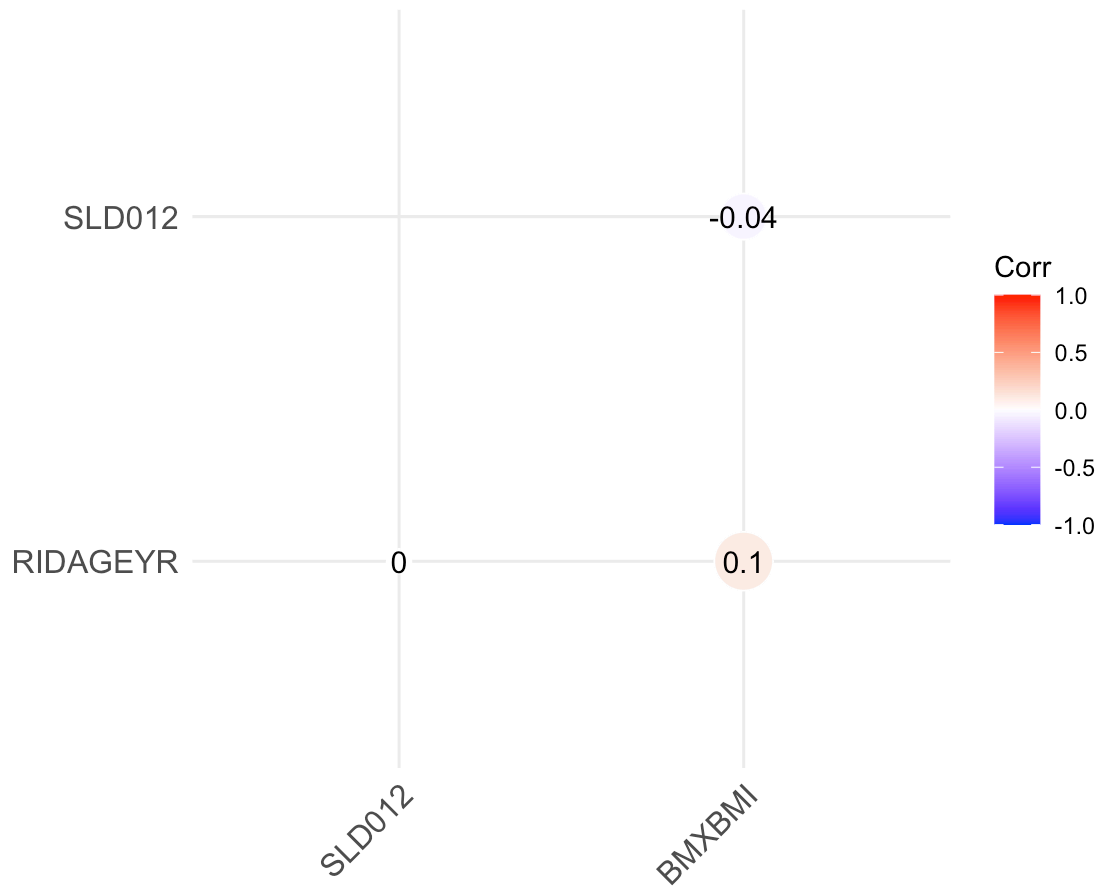
```
## [1] "Correlation Matrix:"
```

```
print(cor_matrix)
```

```
##           RIDAGEYR      SLD012      BMXBMI
## RIDAGEYR  1.0000000000 -0.004756961  0.10081413
## SLD012    -0.004756961  1.0000000000 -0.04462352
## BMXBMI    0.100814131 -0.044623518  1.000000000
```

```
# Plot heatmap of correlation matrix with enhancements
ggcorrplot(cor_matrix,
            method = "circle",      # Circle method for visualization
            type = "lower",        # Only plot the lower triangle of the matrix
            lab = TRUE,            # Show correlation values inside the circles
            lab_size = 4,          # Adjust label size for readability
            colors = c("blue", "white", "red"), # Custom color palette for correlation values
            ggtheme = theme_minimal(), # Minimal theme for a clean look
            outline.col = "white") + # White outline around circles for clarity
ggtitle("Correlation Matrix of Sleep and BMI Variables") + # Title for the plot
theme(plot.title = element_text(size = 14, hjust = 0.5)) # Title size and centering
```

Correlation Matrix of Sleep and BMI Variables



```
options(max.print=15)
# Step 5: Multiple linear regression to see the effect of sleep patterns on BMI
regression_result <- regression_analysis(cleaned_data)
```

```
## [1] "Variance Inflation Factor (VIF) values:"
##           GVIF Df GVIF^(1/(2*Df))
## SLD012    1.021992 1      1.010936
## SLQ030    1.113240 5      1.010785
## SLQ120    1.063247 5      1.006152
## RIDAGEYR  1.057818 1      1.028503
## RIAGENDR  1.029646 1      1.014715
```

## Stp 5) Multi-Linear Regression (simple)

I did a multiple linear regression to check how sleep patterns and other factors like age and gender affect BMI. This helps understand which factors are more important when it comes to BMI.

## Explanation of Output:

### Model Summary:

The regression model shows that sleep frequency (SLQ030), sleep timing (SLQ120), age (RIDAGEYR), and gender (RIAGENDR) all have a significant effect on BMI. For example, if sleep is more frequent (like 5 or more nights a week), BMI tends to be higher. Age and gender also affect BMI, with age slightly increasing BMI and females having a higher BMI than males. The p-value is really low, so the model is statistically significant. But, the R-squared value of 0.1046 means the model only explains about 10.5% of BMI variation, which is not a lot.

### VIF:

The VIF values are all close to 1, which means there's no multicollinearity. In simple terms, the predictor variables aren't overlapping or causing issues in the model.

### Residuals and Plots:

The residuals plot and Q-Q plot help check if the model's assumptions are met. They show if there's any bias or if the data isn't fitting the model properly.

### Key Point:

The regression tells us that sleep habits, age, and gender play a role in BMI. But the model isn't explaining much of the variability in BMI (just 10.5%). Still, the sleep patterns and other factors are significant enough to notice, and we should look into them further.

```
# Print the model summary, residuals plot, and VIF values
print(regression_result$model_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ SLD012 + SLQ030 + SLQ120 + RIDAGEYR + RIAGENDR,
##     data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.755  -4.632  -1.054   3.381  35.838
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                24.766994    0.557014   44.464
## SLD012                   -0.109459    0.056694   -1.931
## SLQ030Rarely - 1-2 nights a week    1.403677    0.245034    5.729
##                                Pr(>|t|)
## (Intercept)                  < 2e-16 ***
## SLD012                       0.05357 .
## SLQ030Rarely - 1-2 nights a week    1.06e-08 ***
## [ reached getOption("max.print") -- omitted 11 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.751 on 5945 degrees of freedom
## Multiple R-squared:  0.1046, Adjusted R-squared:  0.1026
## F-statistic: 53.41 on 13 and 5945 DF, p-value: < 2.2e-16
```

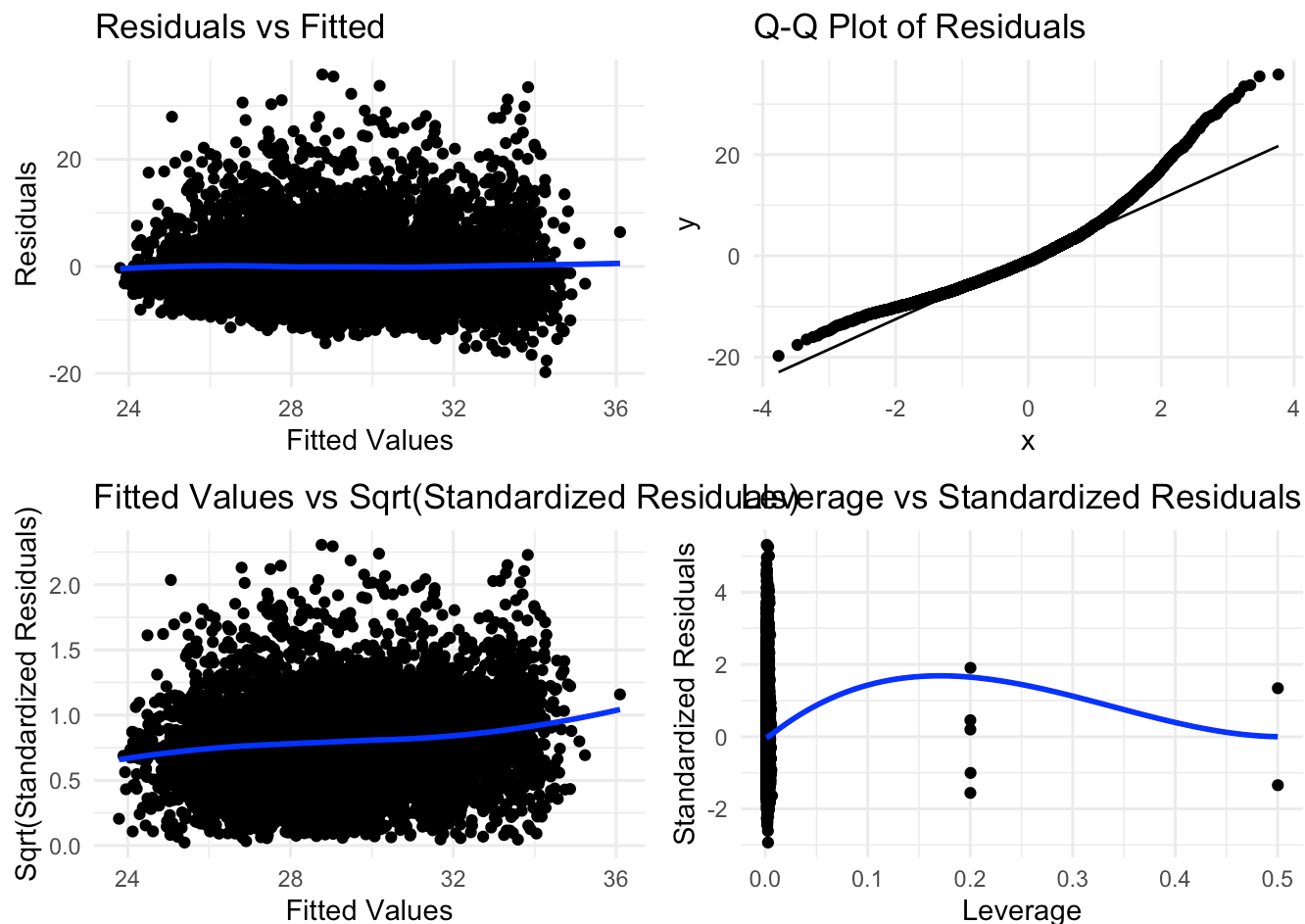
```
print(regression_result$vif_values)
```

```
##              GVIF Df GVIF^(1/(2*Df))
## SLD012      1.021992  1      1.010936
## SLQ030      1.113240  5      1.010785
## SLQ120      1.063247  5      1.006152
## RIDAGEYR    1.057818  1      1.028503
## RIAGENDR    1.029646  1      1.014715
```

```
grid.arrange(regression_result$residuals_plot,
              regression_result$qq_plot,
              regression_result$fitted_vs_sqrt_residuals_plot,
              regression_result$leverage_vs_residuals_plot, ncol = 2)
```

```
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
```





## Step 6) Advanced Multi-Linear Regression Analysis

I did an advanced regression to check for interaction effects between the variables. This helps us understand if the effect of one variable on BMI depends on another variable, like how sleep patterns might change based on age or gender.

### Explanation of Output:

#### Model Summary:

The advanced regression model shows the relationship between BMI and various predictors, considering interaction effects. The multiple R-squared is 0.1483, meaning about 14.8% of the variability in BMI is explained by the model. The adjusted R-squared (0.1178) is a bit lower, showing that after adjusting for the number of predictors, the model still doesn't explain much. The p-value is very small ( $< 2.2e-16$ ), meaning the model as a whole is statistically significant.

#### Residuals and Plots:

The residuals plot, Q-Q plot, and other diagnostic plots help check if the model's assumptions hold true. These plots look at if there's any pattern left in the residuals, which can indicate issues with the model fit.

## Key Point:

The advanced regression model shows significant interaction effects, but like the previous one, it doesn't explain a large portion of BMI variability. The model is statistically significant

```
# Step 6: Advanced Interaction effects tested  
adv_regression_result <- advanced_regression_analysis(cleaned_data)  
  
# Print the model summary, residuals plot, and VIF values  
print(adv_regression_result$model_summary)
```

```

##
## Call:
## lm(formula = BMXBMI ~ SLD012 * SLQ030 * SLQ120 * RIDAGEYR * RIAGENDR,
##     data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.315  -4.451  -0.929   3.321  37.282
##
## Coefficients: (81 not defined because of singularities)
##
Estimate
## (Intercept)
3.028e+01
## SLD012
-8.910e-01
## SLQ030Rarely - 1-2 nights a week
-1.059e+01
##
Std. Error
## (Intercept)
5.970e+00
## SLD012
7.068e-01
## SLQ030Rarely - 1-2 nights a week
1.051e+01
##
t value
## (Intercept)
5.072
## SLD012
-1.261
## SLQ030Rarely - 1-2 nights a week
-1.008
##
Pr(>|t|)
## (Intercept)
4.06e-07
## SLD012
0.2075
## SLQ030Rarely - 1-2 nights a week
0.3137
##
## (Intercept)
***
## SLD012
## SLQ030Rarely - 1-2 nights a week
## [ reached getOption("max.print") -- omitted 285 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.693 on 5752 degrees of freedom

```

```
## Multiple R-squared:  0.1483, Adjusted R-squared:  0.1178  
## F-statistic: 4.861 on 206 and 5752 DF,  p-value: < 2.2e-16
```

```
grid.arrange(adv_regression_result$residuals_plot,  
              adv_regression_result$qq_plot,  
              adv_regression_result$fitted_vs_sqrt_residuals_plot,  
              adv_regression_result$leverage_vs_residuals_plot, ncol = 2)
```

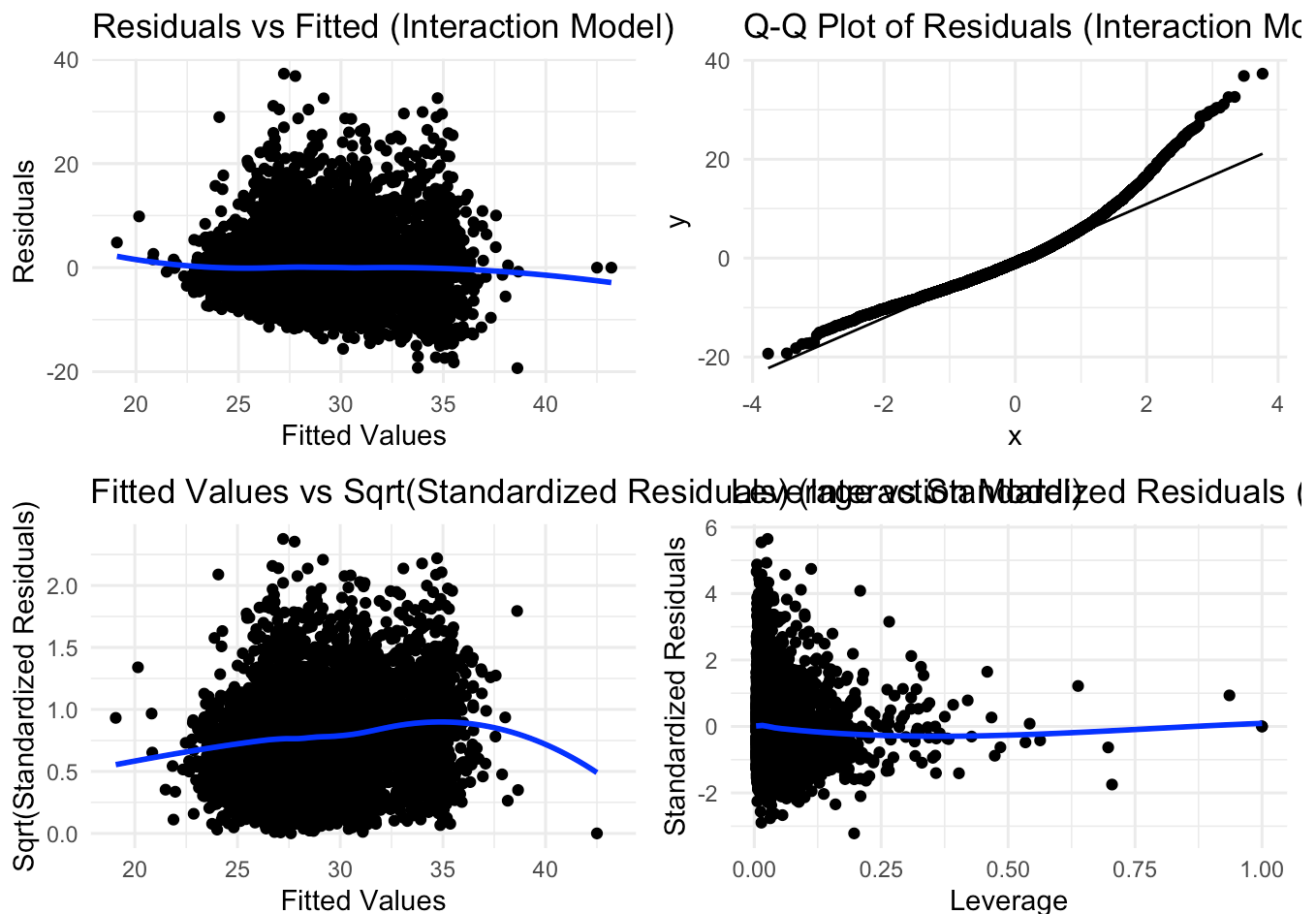
```
## `geom_smooth()` using formula = 'y ~ x'  
## `geom_smooth()` using formula = 'y ~ x'
```

```
## Warning: Removed 5 rows containing non-finite outside the scale range  
## (`stat_smooth()`).
```

```
## Warning: Removed 5 rows containing missing values or values outside the scale range  
## (`geom_point()`).
```

```
## `geom_smooth()` using formula = 'y ~ x'
```

```
## Warning: Removed 5 rows containing non-finite outside the scale range (`stat_smooth()  
)`.  
## Removed 5 rows containing missing values or values outside the scale range  
## (`geom_point()`).
```



## Step 7) Subgroup Analysis

In this step, I did a subgroup analysis to see how demographic factors (like sleep duration, snoring, and sleepiness) affect the relationship between sleep patterns and BMI. The idea is to break down the data by categories and see how these factors interact with BMI.

## Explanation of Output:

### Regression Line Plot:

The plot `plot_bmi_sleep_age` shows how BMI changes with sleep duration and age. This helps visualize the relationship between sleep and BMI in different age groups.

## Key Point:

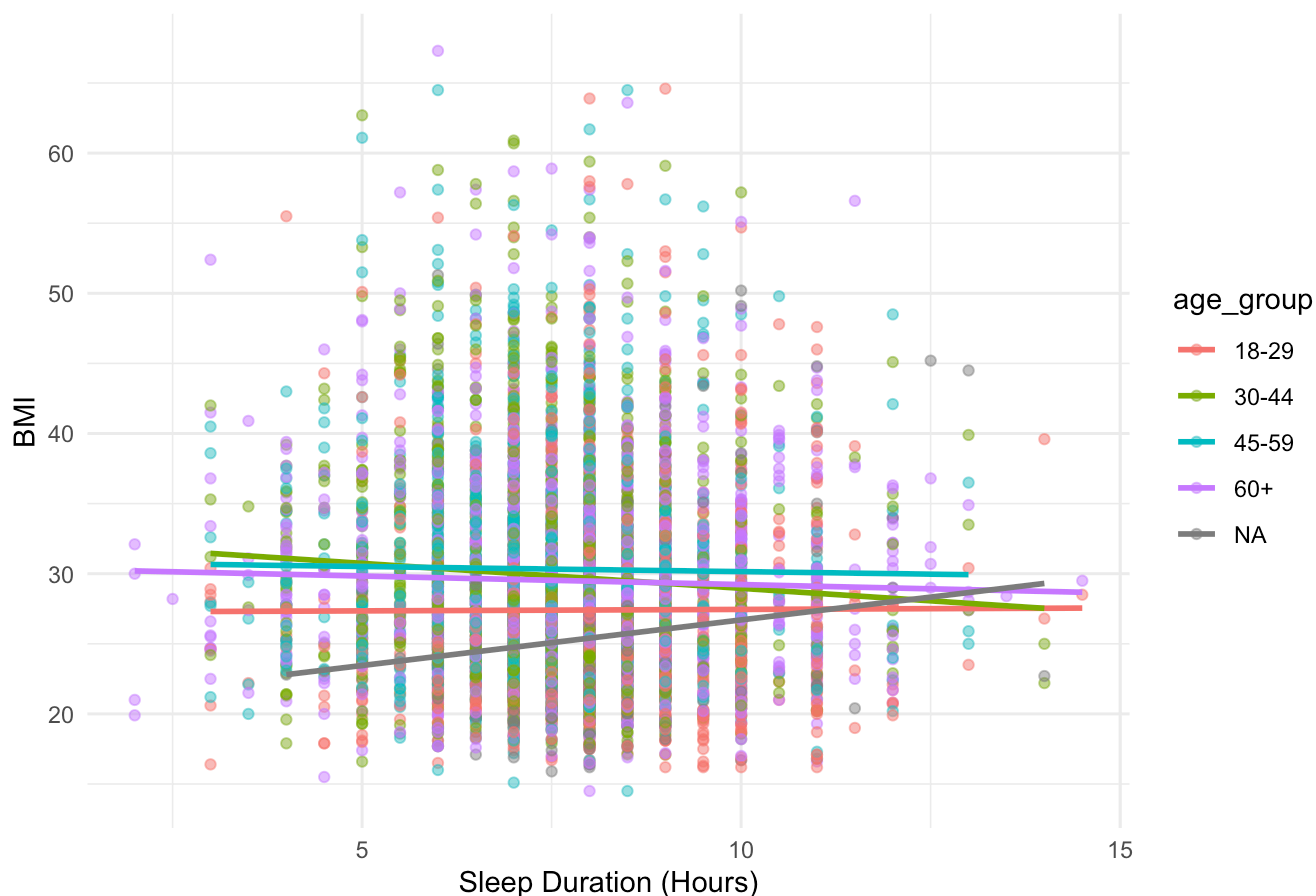
In this step, demographic factors like snoring and sleepiness have a more significant effect on BMI than sleep duration. Although these variables show statistically significant relationships, the amount of variability in BMI they explain is still quite small. There may be other factors not included in the model that are affecting BMI more strongly. This analysis suggests that while sleep patterns do influence BMI, the impact is modest and may need further exploration with additional variables or different analysis methods.

```
# Step 7: Subgroup analysis to see how demographic factors influence the relationship
#result <- subgroup_analysis(cleaned_data)
# Call the regression line plots function
# Call the regression line plots function
regression_results <- regression_line_plots(cleaned_data)

# Print the plots
print(regression_results$plots$plot_bmi_sleep_age)
```

```
## `geom_smooth()` using formula = 'y ~ x'
```

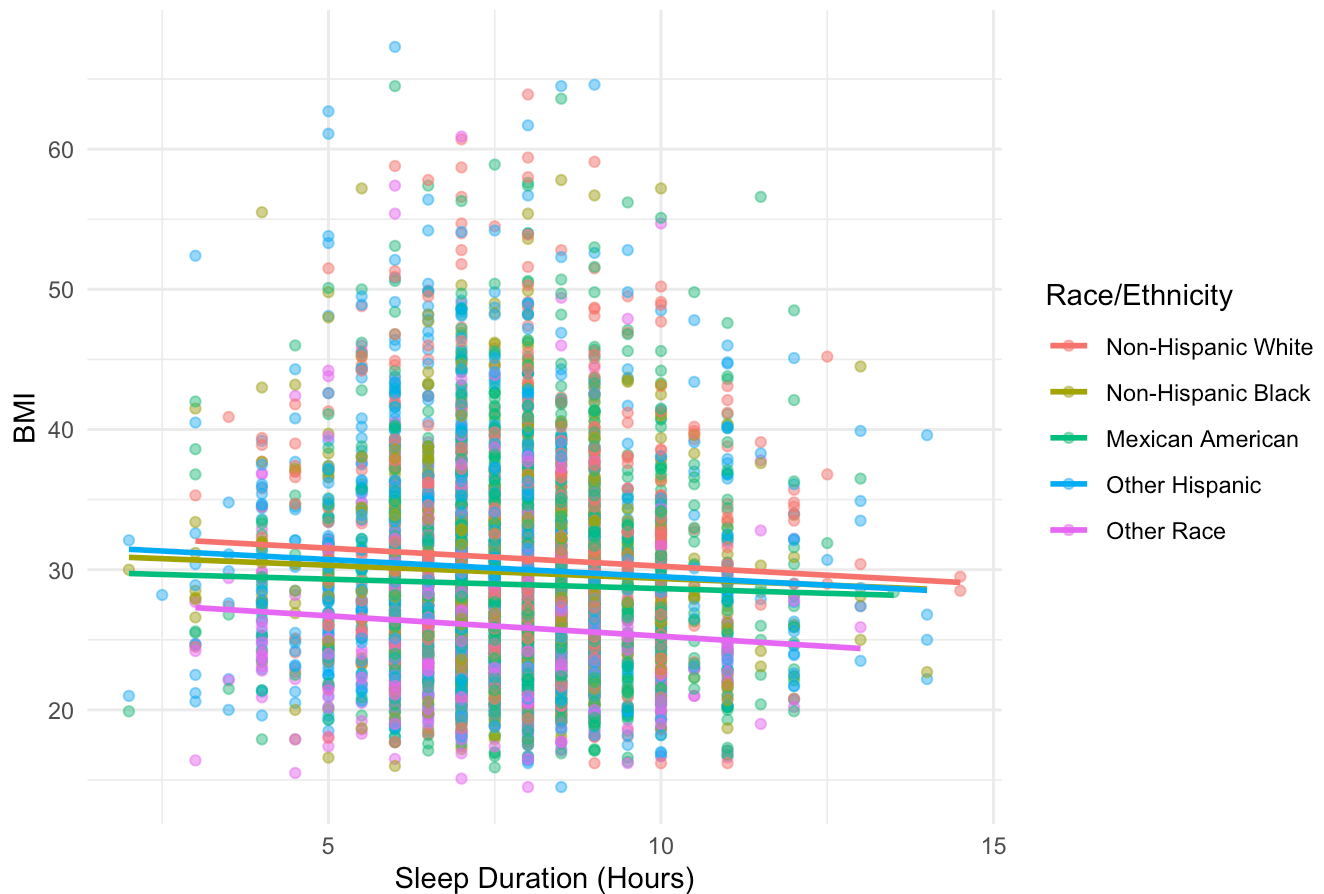
Regression of BMI vs Sleep Duration by Age Group



```
print(regression_results$plots$plot_bmi_sleep_race)
```

```
## `geom_smooth()` using formula = 'y ~ x'
```

## Regression of BMI vs Sleep Duration by Race/Ethnicity



```
# Print the summaries
print(regression_results$summaries$age_group_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ SLD012 * age_group, data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.761  -5.071  -1.170   3.532  37.585
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    27.24414     1.11729   24.384 < 2e-16 ***
## SLD012          0.02102     0.13512    0.156 0.876413
## age_group30-44    5.27774     1.48461    3.555 0.000381 ***
## [ reached getOption("max.print") -- omitted 5 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.027 on 5627 degrees of freedom
## (324 observations deleted due to missingness)
## Multiple R-squared:  0.02171,    Adjusted R-squared:  0.02049
## F-statistic: 17.84 on 7 and 5627 DF,  p-value: < 2.2e-16
```

```
print(regression_results$summaries$ethnicity_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ SLD012 * factor(RIDRETH1), data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.372  -4.778  -1.167   3.552  36.818
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    32.82461     1.12289   29.232 < 2e-16
## SLD012         -0.25721     0.13975   -1.841  0.06573
## factor(RIDRETH1)Non-Hispanic Black -1.56444     1.62848   -0.961  0.33675
##
## (Intercept)          ***
## SLD012                .
## factor(RIDRETH1)Non-Hispanic Black
## [ reached getOption("max.print") -- omitted 7 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.954 on 5949 degrees of freedom
## Multiple R-squared:  0.04929,    Adjusted R-squared:  0.04785
## F-statistic: 34.27 on 9 and 5949 DF,  p-value: < 2.2e-16
```

```
# Call the boxplots function
boxplot_results <- subgroup_boxplots(cleaned_data)

# Print the plots
```

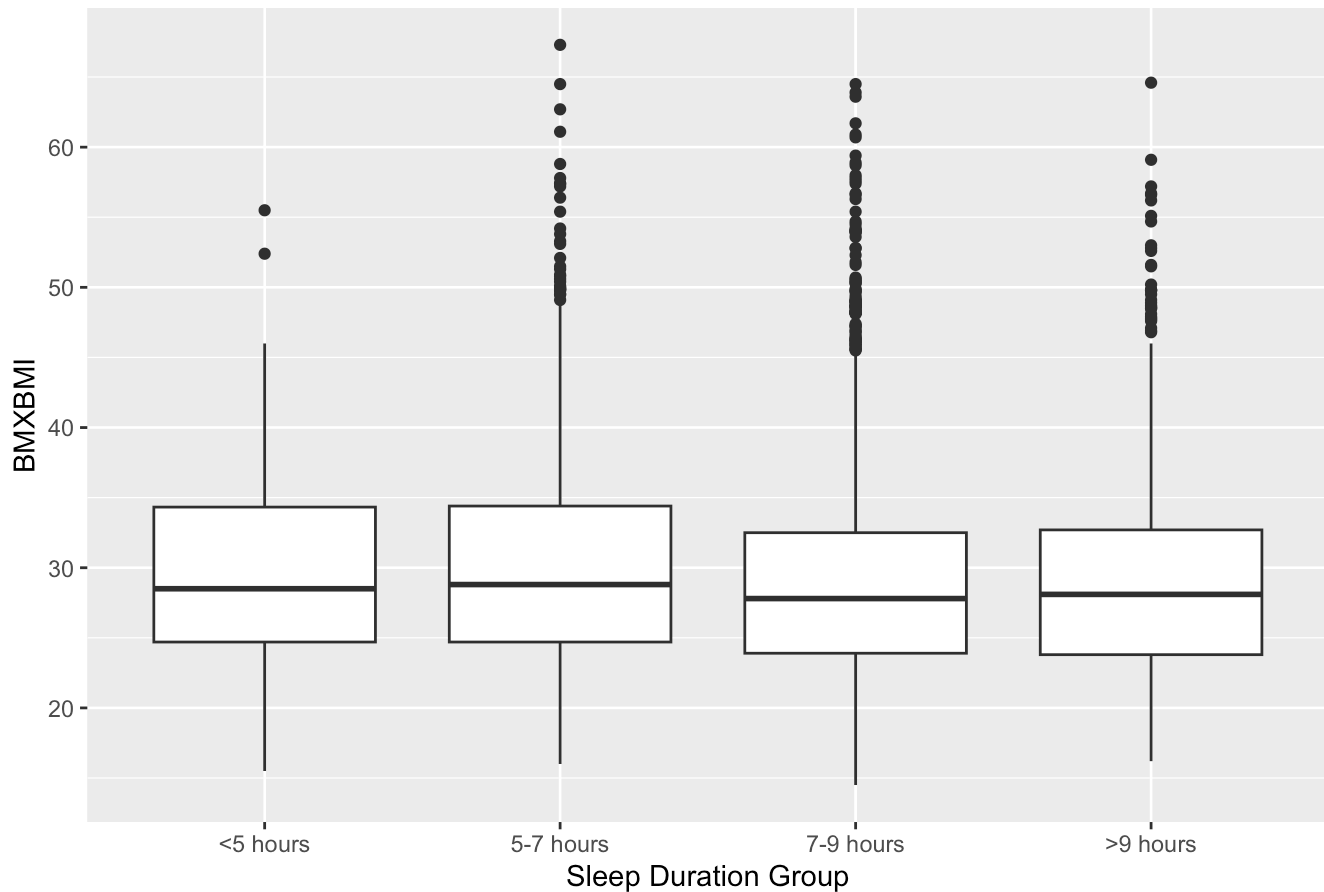
## Boxplot Summaries for Sleep Duration:

The regression model shows that sleep duration (grouped by hours) has a small effect on BMI. The p-value for the sleep duration groups (5-7 hours, 7-9 hours) is higher than 0.05, indicating that these groups do not significantly affect BMI. The multiple R-squared value is very low (0.004628), suggesting that sleep duration doesn't explain much of the variation in BMI. The F-statistic and p-value show that the model is significant, but the actual effect size is very small.

```
print(boxplot_results$plots$sleep_duration)
```



BMI by Sleep Duration

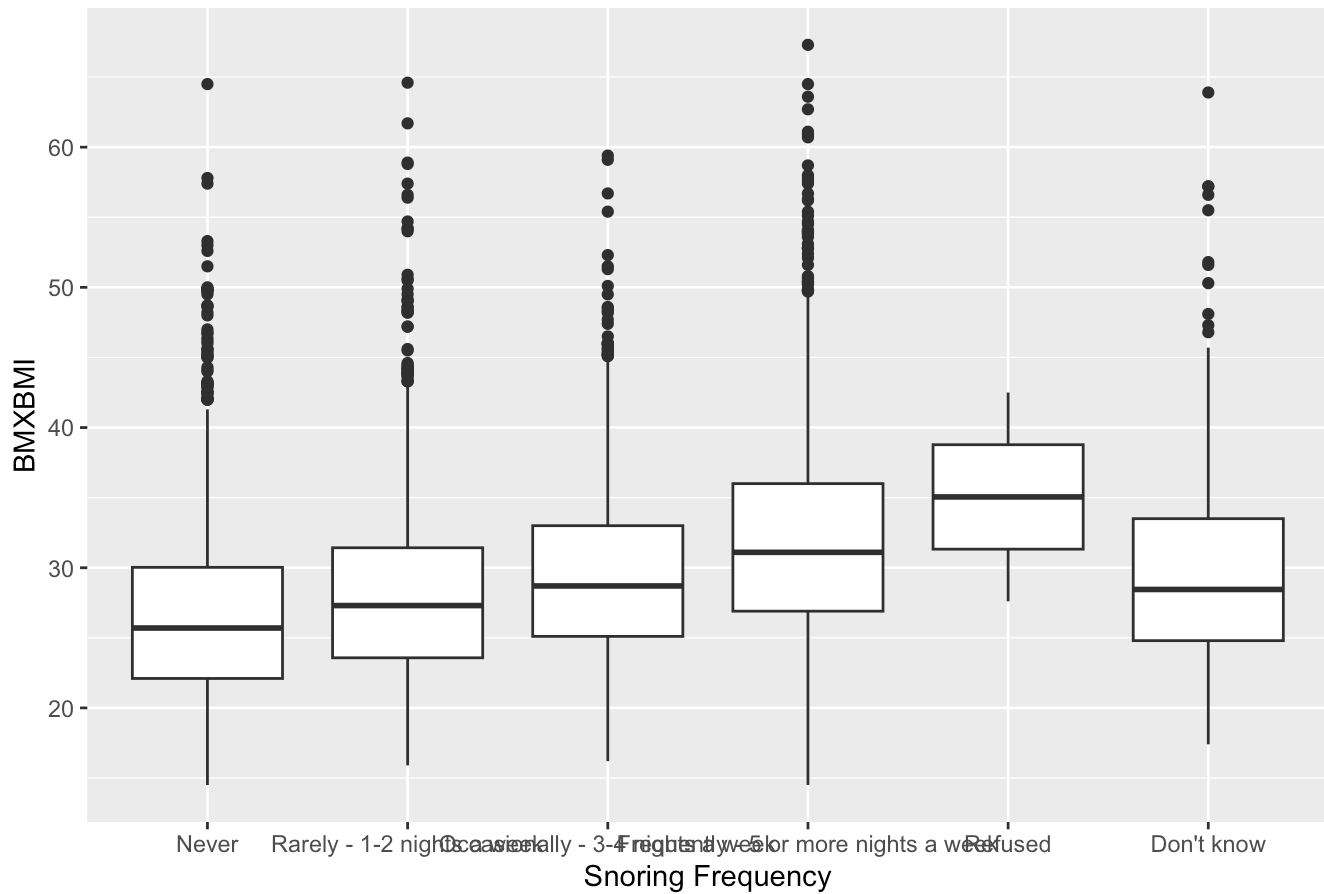


## Boxplot Summaries for Snoring:

The analysis of snoring (grouped by frequency of snoring nights) shows that snoring is more significant. The estimates for groups like “Rarely - 1-2 nights a week” and “Occasionally - 3-4 nights a week” have p-values well below 0.05, indicating these categories significantly affect BMI. The R-squared (0.07946) is slightly better here, suggesting snoring explains a bit more of the variability in BMI.

```
print(boxplot_results$plots$snoring)
```

BMI by Snoring Frequency

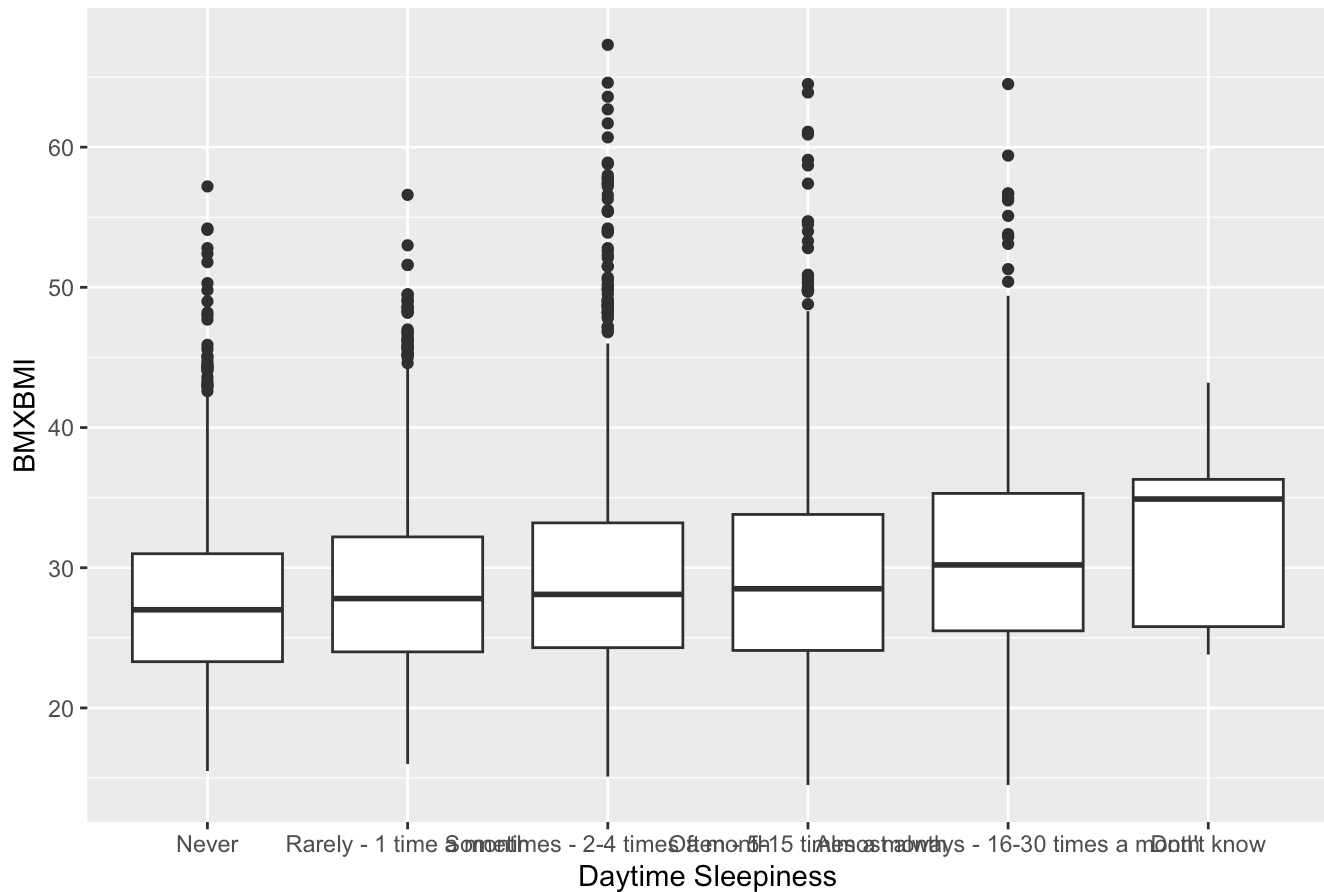


## Boxplot Summaries for Sleepiness:

The sleepiness groups show a similar pattern. “Rarely - 1 time a month” and “Sometimes - 2-4 times a month” have significant effects on BMI, with very low p-values. The R-squared (0.01565) is still low, but the p-values indicate a statistically significant relationship between sleepiness and BMI.

```
print(boxplot_results$plots$sleepiness)
```

## BMI by Daytime Sleepiness



```
# Print the summaries
print(boxplot_results$summaries$sleep_duration_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ sl_sleep_group, data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14.380  -5.080  -1.080   3.723  37.168
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      29.5672     0.5301  55.781  <2e-16 ***
## sl_sleep_group5-7 hours    0.5648     0.5706   0.990    0.322
## sl_sleep_group7-9 hours   -0.6877     0.5447  -1.263    0.207
## [ reached getOption("max.print") -- omitted 1 row ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.111 on 5955 degrees of freedom
## Multiple R-squared:  0.004628,    Adjusted R-squared:  0.004126
## F-statistic: 9.229 on 3 and 5955 DF,  p-value: 4.353e-06
```

```
print(boxplot_results$summaries$snoring_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ factor(SLQ030), data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17.706  -4.806  -1.074   3.460  37.660
##
## Coefficients:
##                                     Estimate Std. Error t value
## (Intercept)                       26.8397     0.1665 161.213
## factor(SLQ030)Rarely - 1-2 nights a week      1.4340     0.2471   5.804
## factor(SLQ030)Occasionally - 3-4 nights a week  2.8291     0.2725  10.381
##                                     Pr(>|t|)
## (Intercept)                       < 2e-16 ***
## factor(SLQ030)Rarely - 1-2 nights a week      6.81e-09 ***
## factor(SLQ030)Occasionally - 3-4 nights a week < 2e-16 ***
## [ reached getOption("max.print") -- omitted 3 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.84 on 5953 degrees of freedom
## Multiple R-squared:  0.07946,    Adjusted R-squared:  0.07869
## F-statistic: 102.8 on 5 and 5953 DF,  p-value: < 2.2e-16
```

```
print(boxplot_results$summaries$sleepiness_summary)
```

```
##
## Call:
## lm(formula = BMXBMI ~ factor(SLQ120), data = cleaned_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16.560  -5.066  -1.066   3.634  37.791
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                   27.7698     0.2151 129.083
## factor(SLQ120)Rarely - 1 time a month      0.9197     0.2873   3.201
## factor(SLQ120)Sometimes - 2-4 times a month  1.7389     0.2686   6.473
##                                Pr(>|t|)
## (Intercept)                   < 2e-16 ***
## factor(SLQ120)Rarely - 1 time a month      0.00138 **
## factor(SLQ120)Sometimes - 2-4 times a month 1.04e-10 ***
## [ reached getOption("max.print") -- omitted 3 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.073 on 5953 degrees of freedom
## Multiple R-squared:  0.01565,    Adjusted R-squared:  0.01482
## F-statistic: 18.93 on 5 and 5953 DF,  p-value: < 2.2e-16
```

## Summary Analysis of this study

This analysis explores the relationship between sleep patterns—such as sleep duration, snoring frequency, and daytime sleepiness—and key health measures like Body Mass Index (BMI) in adults. The goal was to uncover connections between these sleep behaviors and BMI, as well as how demographic factors such as age, gender, and ethnicity might affect these relationships.

### Primary Study Aim

The main aim of this study was to understand how sleep patterns correlate with BMI. Through our analysis, I examined variables like sleep duration, snoring, and daytime sleepiness in relation to BMI. Several steps were involved in this process:

**Data Loading and Cleaning:** I first loaded and cleaned datasets related to demographic information, sleep quality, and BMI. This included merging data from the DEMO, SLQ, and BMX files to ensure a clean dataset for analysis.

**Descriptive Statistics for BMI:** I computed and reviewed basic statistics for BMI in the population, showing the mean, median, and range for BMI.

**Correlation Analysis:** A correlation analysis revealed moderate relationships between sleep patterns (such as sleep duration, snoring frequency, and daytime sleepiness) and BMI. The correlation matrix was visualized using a heatmap, providing useful insights.

**Multiple Linear Regression:** A regression model assessed how sleep patterns affect BMI. The results indicated that snoring and sleep duration had moderate associations with BMI, while daytime sleepiness did not significantly influence BMI. Although the regression model explained only a small portion of the variability in BMI

(R-squared of 0.004), the insights remain valuable.

## Secondary Study Aim

The secondary objective was to explore how demographic factors like age, gender, and ethnicity influence the relationships between sleep behaviors and health outcomes. Subgroup analysis revealed that these demographic factors interact with the primary sleep metrics.

Subgroup Analysis: Regression line plots showed that age and ethnicity influence the relationship between sleep patterns and BMI. Significant variations were observed across different demographic groups.

Boxplots and Summaries: Boxplots visualized differences in sleep duration, snoring frequency, and sleepiness across demographic groups, revealing varying patterns that affect BMI differently.

## Statistical Summary of Findings

Correlation Analysis: The correlation matrix showed modest relationships between sleep duration and BMI, with longer sleep durations moderate associated with lower BMI. Snoring and daytime sleepiness followed similar trends.

Regression Results: The regression analysis showed that sleep duration and snoring had associations with BMI. The model's Adjusted R-squared value indicates that other factors may play a more significant role in influencing BMI.

Subgroup Analysis: Age and ethnicity moderately influenced the relationship between sleep patterns and BMI, The variations were borderline large enough to draw this conclusion.

## Hypotheses Review

Primary Null Hypothesis (H0): The hypothesis that no significant relationship exists between sleep patterns and BMI was rejected. While relationships were identified, they were moderate, suggesting more complex models may be required for a fuller understanding.

Secondary Null Hypothesis (H0): The hypothesis that demographic factors do not influence the relationship between sleep and BMI was rejected. Age and ethnicity showed moderate influence.

Primary Alternative Hypothesis (H1): The hypothesis that shorter sleep durations and poor sleep quality are associated with higher BMI was supported with weak to moderate associations.

Secondary Alternative Hypothesis (H1): The hypothesis that demographic factors influence the relationship between sleep and BMI was supported. Age and ethnicity demonstrated moderating effects on sleep-BMI relationships, though further investigation is needed for deeper insights.

## Conclusion

The analysis demonstrates statistically significant relationships between sleep patterns and BMI, although with weak to moderate associations. The models highlight important trends, such as correlations between BMI and both sleep duration and snoring frequency. Additionally, demographic factors like age and ethnicity appear to influence these relationships.

## Next step:

Strengthening the Findings However, stronger associations or confidence could be achieved by considering additional factors such as:

Lifestyle Habits (diet, exercise, and stress levels) Mental Health Factors (anxiety, depression) Environmental Influences (sleep environment, work schedules) Genetic Predispositions Incorporating these factors in future research could provide a clearer and more comprehensive understanding of the relationships between sleep, demographic factors, and health outcomes like BMI.