# Sleep Efficiency Prediction Using Linear Regression

Weiqing Xu[a], Haolin Cen[a], and Yuzhuo Liu[a]

[a]Shanghai Jiao Tong University, 800 Dongchuan Rd., Shanghai, China

## ABSTRACT

There are various of variables which can be used to evaluate the quality of individuals' sleep. Among those, sleep efficiency (SE) is most commonly used. It has garnered much attention in many studies and shows a significance to people's physical and mental health. Moreover, there are several stages of sleep such as deep sleep, which is also an essential factor. In this report, we will use linear regression together with various of techniques to show you the correlations between SE, different stages of sleep, and individuals' health habits. We will show our motivation and description first, construct three models step by step, evaluate the performance, and draw the conclusion.

**Keywords:** linear regression, data analysis, regression analysis, R, sleep efficiency, deep sleep

## 1. INTRODUCTION

### 1.1 Background

Sleep efficiency (SE) is the most commonly used measurement to evaluate the quality of the sleep [1]. Previous researches have proved that SE is more related to better sleeping compared to other factors [2]. When it comes to measuring SE, deep sleep percentage is widely used and regarded as one of the most important measuring factor [3].

The monitoring of SE has garnered significant attention in medicine and healthcare due to its association with various health conditions. Studies have linked low SE to snoring, asthma, physical health-related quality of life, interstitial lung disease, cognitive impairment, spinal cord injury, attention deficit hyperactivity disorder, mild cognitive impairment, and Alzheimer's disease [4-10]. These studies showed the significance of sleeping efficiency to people's physical and mental health.

As a result, there is an urge for us to find out more factors that may influence sleeping efficiency. Our study aims to find out more relationships of predictors and SE, as well as if deep sleep percentage is a good measurement of SE. In our study, we analyze a data set which contains numbers of factors such as age, gender, awakenings, caffeine consumption, alcohol consumption, smoking status and exercise frequency. Firstly, we construct models to reveal the correlation between sleeping efficiency and other factors. Then, an analysis of deep sleep percentage is taken to evaluate whether it is a good measure of SE and find out how it is determined. Finally, we can find out the relationship between SE and deep sleep percentage based on our analyses.

### 1.2 Data Set

In this study, we use a data set which contains reasonable number of cases (around 500 samples). Below is a part of the data for reference and the data is new enough for statistical analysis.

In this data set, each test subject is identified by "ID". Personal information such as "Age" and "Gender" is also recorded. The "Bedtime" and "Wakeup time" indicate when each subject goes to bed and wakes up each day, and the "Sleep duration" records the total amount of time each subject slept in hours. "Sleep Efficiency" is calculated as the ratio between actual sleep time and time spent in bed. The "REM sleep percentage", "Deep sleep percentage", and "Light sleep percentage" indicate the amount of time spent in each stage of sleep. The "Awakenings" records the number of times each subject wakes up during the night. Additionally, the data set includes information about each subject's caffeine and alcohol consumption (ml) in the 24 hours prior to bedtime, their smoking status (Yes or No), and their exercise frequency (times per week).

| ID | Age | Gender | Bedtime | Wakeup time | Sleep duration | Sleep efficiency | REM sleep pe | Deep sleep p | Light sleep peri | Awakenings | Caffeine consumption | Alcohol consumption | Smoking status | Exercise frequency |
|----|-----|--------|---------|-------------|----------------|------------------|--------------|--------------|------------------|------------|----------------------|---------------------|----------------|--------------------|
| 1 | 65 | Female | 2021/3/6 1:00 | 2021/3/6 7:00 | 6 | 0.88 | 18 | 70 | 12 | 0 | 0 | 0 | Yes | 3 |
| 2 | 69 | Male | 2021/12/5 2:00 | 2021/12/5 9:00 | 7 | 0.66 | 19 | 28 | 53 | 3 | 0 | 3 | Yes | 3 |
| 3 | 40 | Female | 2021/5/25 21:30 | 2021/5/25 5:30 | 8 | 0.89 | 20 | 70 | 10 | 1 | 0 | 0 | No | 3 |
| 4 | 40 | Female | 2021/11/3 2:30 | 2021/11/3 8:30 | 6 | 0.51 | 23 | 25 | 52 | 3 | 50 | 5 | Yes | 1 |

# 2. METHODS & RESULTS

## 2.1 Modeling Sleeping Efficiency

### 2.1.1 Overview

In this part, we will develop a model studying the correlation between sleep efficiency and other lifestyle variables including:
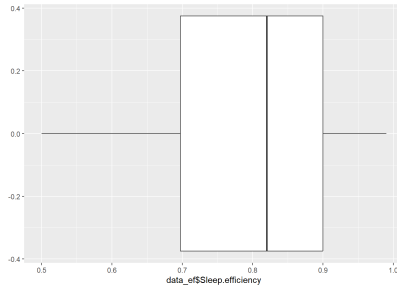
- Sleeping time and the number of awakenings, which are variables directly related to sleep efficiency,
- Some personal information such as gender and age,
- Individual's daily dietary, smoking, exercise, and other Health Habits.

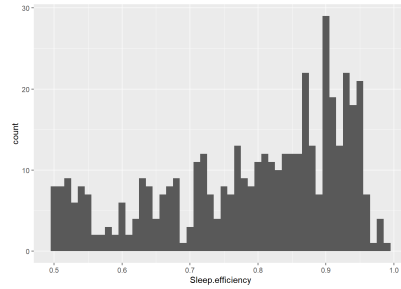### 2.1.2 Data preprocessing and exploration

Firstly, we analyzed the statistical data of sleep efficiency (table 1) and provide a bar plot (figure 1 (a) and (b)). From the histogram we can see that the distribution is left skewed, which aligns with the statistics table that mean is less than median.

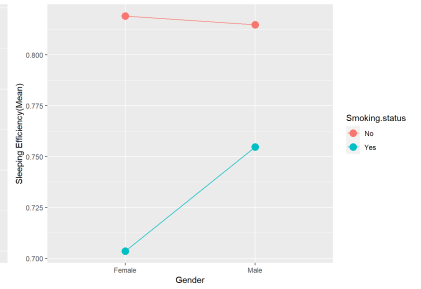| mean | median | variance | standard deviation | min | max |
|------|--------|----------|--------------------|-----|-----|
| 0.7889159 | 0.82 | 0.01828907 | 0.1352371 | 0.5 | 0.99 |

Table 1. Statistics for sleep efficiency



(a) Box Plot   (b) Histogram   (c) Means of sleeping efficiency

Figure 1. Box plot and histogram of sleeping efficiency

Then, we notice that there are some dummy variables in these predictors such as gender and smoking status, and these two aspects divide the data into four categories. We can observe from figure 1 (c) that the means are different. Additionally, we use box plot and histogram for each category. See figure 2.

Last, we notice that there are some data missing. The strategy we use is to just omitting the corresponding entry, since adding data generated either randomly or by statistics of other data will lead to data contamination, which may violate the assumptions of linear models, and adversely affects our modeling process.

### 2.1.3 Base Model

After basic data preprocessing and exploration, we have a base idea and build our base model **Sleep.efficiency = Sleep.duration + Age + Gender + Awakenings + Caffeine.consumption + Alcohol.consumption + Smoking.status + Exercise.frequency**.
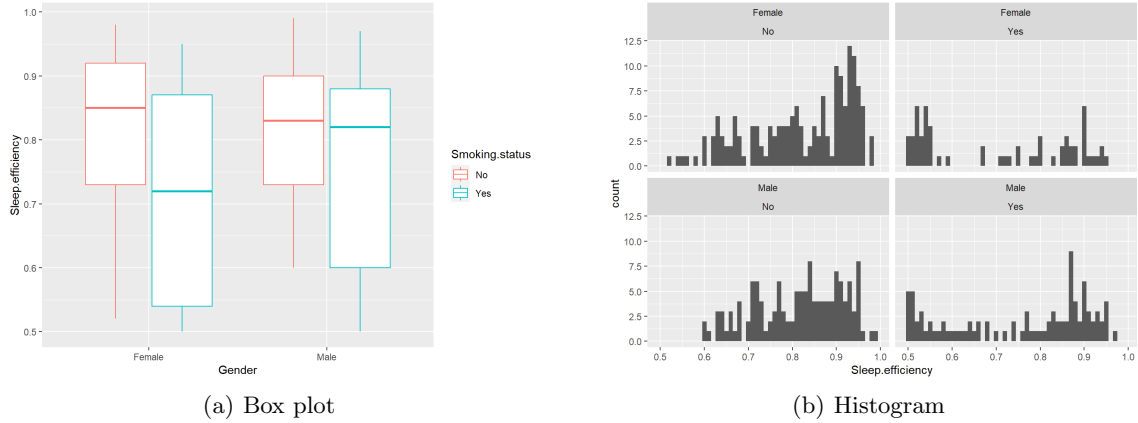
(a) Box plot

(b) Histogram

Figure 2. Box plot and histogram of sleeping efficiency



(a) Box Cox for base model

(b) Box Cox for the model with transformation
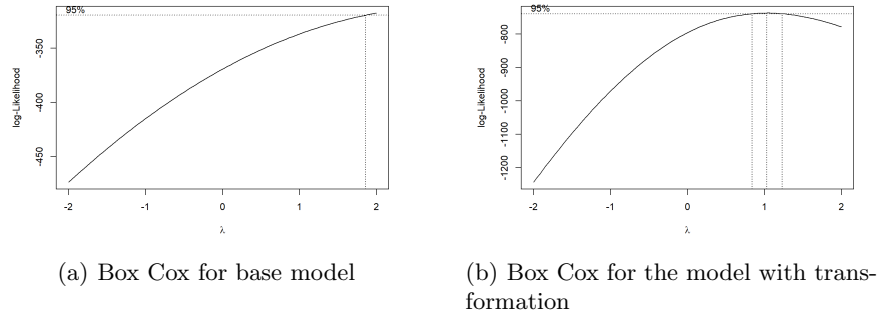
Figure 3. Box plot and histogram of sleeping efficiency

```
                         Estimate    Std. Error  t value  Pr(>|t|)
(Intercept)             0.8565715    0.0470398   18.209   < 2e−16  ***
Sleep.duration         −0.0028873    0.0055022   −0.525   0.60006
Age                     0.0013656    0.0003785    3.608   0.00035  ***
GenderMale              0.0067363    0.0108019    0.624   0.53325
Awakenings             −0.0475818    0.0038088  −12.493   < 2e−16  ***
Caffeine.consumption    0.0001719    0.0001763    0.975   0.33027
Alcohol.consumption    −0.0238733    0.0030974   −7.708   1.14e−13  ***
Smoking.statusYes      −0.0788579    0.0103676   −7.606   2.25e−13  ***
Exercise.frequency      0.0128543    0.0035628    3.608   0.00035  ***
Residual standard error: 0.09506 on 379 degrees of freedom
Multiple R−squared:  0.5195, Adjusted R−squared:  0.5094
F−statistic: 51.22 on 8 and 379 DF,  p−value: < 2.2e−16
```

From the output above we can find that the model pass the F-test. As for each predictors, only age, awakenings, alcohol, smoking, and exercise frequency pass t-test.

### 2.1.4 Transformation of response variable

Next, we consider using Box-Cox method to automatically select the "best" power $\lambda$ that make the residuals of the model closest to normal and constant variability, plotted in Figure 3 (a).

We can see the maximum value of log-Likelihood is achieved when the power is approximately 3. After that, we change the model to **Sleep.efficiency$^3$ = Sleep.duration + Age + Gender + Awakenings + Caffeine.consumption + Alcohol.consumption + Smoking.status + Exercise.frequency**

```
Residual standard error: 0.1607 on 379 degrees of freedom
Multiple R-squared:  0.5412, Adjusted R-squared:  0.5315
F-statistic: 55.88 on 8 and 379 DF,  p-value: < 2.2e-16
```

In the above R output, although the residual standard error increases due to the influence of the cubic term, the coefficient of determination (R-squared) is improved, which means the model can explain more proportion of variance in sleeping efficiency. The result of Box Cox method on the new model is in Figure 3 (b).

### 2.1.5 Feature Selection

As stated above, there are some predictors which does not pass t-test. Since we just put all seemingly factors as predictors, we should eliminate some insignificant varibles. The technique we use here is backward elimination (BE) and we do not consider interactions. After the BE process stops, it gives the model **Sleeping.efficiency**$^3$ **= Age + Awakenings + Alcohol.consumption + Smoking.status + Exercise.frequency**

| (Intercept) | Age | Awakenings | Alcohol.consumption | SmoYes | Exercise.frequency |
|---|---|---|---|---|---|
| 0.63671 | 0.00209 | -0.09121 | -0.03802 | -0.11597 | 0.02418 |

Table 2. model parameters

We can see from the table 2 that parameter result that the remaining significant predictors are the same as those which pass t-test. Then, we analyze this model and compare it with the full model.

```
## Residual standard error: 0.1604 on 382 degrees of freedom
## Multiple R-squared:  0.5394, Adjusted R-squared:  0.5333
## F-statistic: 89.46 on 5 and 382 DF,  p-value: < 2.2e-16
##
## Model 1: Sleep.efficiency^3 ~ Age + Awakenings + Alcohol.consumption +
##     Smo + Exercise.frequency
## Model 2: Sleep.efficiency^3 ~ Sleep.duration + Age + Gen + Awakenings +
##     Caffeine.consumption + Alcohol.consumption + Smo + Exercise.frequency
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1    382 9.8225
## 2    379 9.7839  3  0.038579 0.4981 0.6838
```

From the above R output we can see that the adjusted R-squared slightly increases. For the anova table, it performs a F-test, where the null hypothesis chooses the reduced model, and the results tells that we cannot reject null hypothesis. Consequently, the reduced model is better,

For interactions, as we showed in figure 2, there might be some interactions between gender and smoking status. So we perform a test and the result is listed below:

```
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1    382 9.8225
## 2    380 9.7926  2  0.029958 0.5813 0.5597
```

We can see that although we guess there is an interaction between dummy variables, adding the interaction cannot pass F-test compared with the model without it.

### 2.1.6 Statistical inference and parameter interpretation

Now we can use this model to do some predictions, but first, we should take a look at its confidence interval to see the accuracy of this model. The R output for 95% confidence interval for parameters is listed in figure 4 (a).

As for prediction, we take a 21 years old gentleman as an example. He never smokes or drinks, sleeps without awakenings during the night, and exercises once every week. The confidence interval and prediction interval for his sleeping efficiency is listed in figure 4 (b).
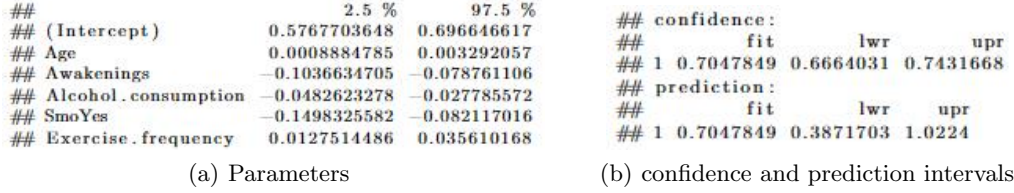
```
##                          2.5 %        97.5 %
## (Intercept)          0.5767703648   0.696646617
## Age                  0.0008884785   0.003292057
## Awakenings          -0.1036634705  -0.078761106
## Alcohol.consumption -0.0482623278  -0.027785572
## SmoYes              -0.1498325582  -0.082117016
## Exercise.frequency   0.0127514486   0.035610168
```

(a) Parameters

```
## confidence:
##           fit        lwr        upr
## 1  0.7047849  0.6664031  0.7431668
## prediction:
##           fit        lwr        upr
## 1  0.7047849  0.3871703  1.0224
```

(b) confidence and prediction intervals

Figure 4. Confidential Intervals

Notice that the sleeping efficiency cannot exceeds 1 due to its definition. Thus, these results mean that, the average sleeping efficiency for him (or persons have same corresponding characteristics, generally) is 0.705, ranging from 0.666 to 0.743 with 95% confidence, and the actual sleeping efficiency is predicted between 0.387 and 1. We can see that his sleeping efficiency exhibits significant uncertainty, with a considerable range of fluctuations.

### 2.1.7 Regression Diagnostics

Now, we cannot be sure that this is our final stage for this modelling process, since we need to test whether regression assumptions are violated and the quality of the data of predictors. We use R method to test collinearity (using *vif* in R) and outliers. We find that there is no raw or studentized or externally studentized residuals that are above 3, which shows that there is no outliers. Table 3 shows that there is no obvious linearity relationship between indicators.

| Age | Awakenings | Alcohol.consumption | Smo | Exercise.frequency |
|---|---|---|---|---|
| 1.010116 | 1.109717 | 1.061496 | 1.008010 | 1.066041 |

Table 3. vif result for testing collinearity

### 2.1.8 Model Evaluation

Finally, we can end up this modeling part by model evaluation. The root mean squared error, R-squared, F-statistics, as well as residual vs. fitted values plot are shown in figure 5.

```
## Residual standard error: 0.1604 on 382 degrees of freedom
## Multiple R-squared:  0.5394, Adjusted R-squared:  0.5333
## F-statistic: 89.46 on 5 and 382 DF,  p-value: < 2.2e-16
```

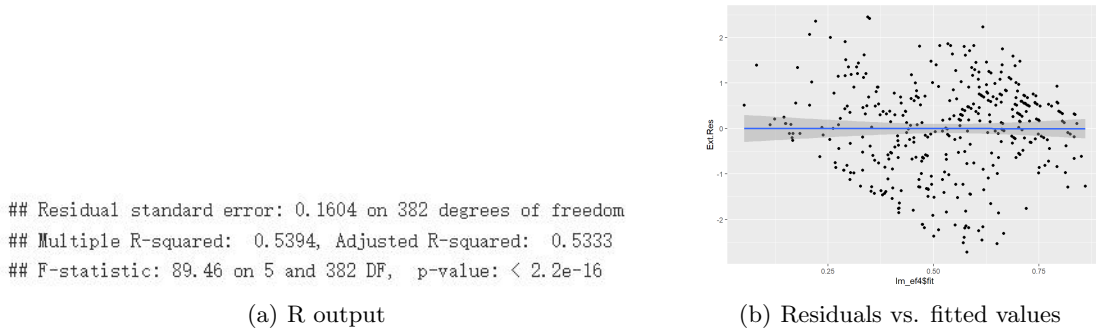(a) R output



(b) Residuals vs. fitted values

Figure 5. Model evaluation

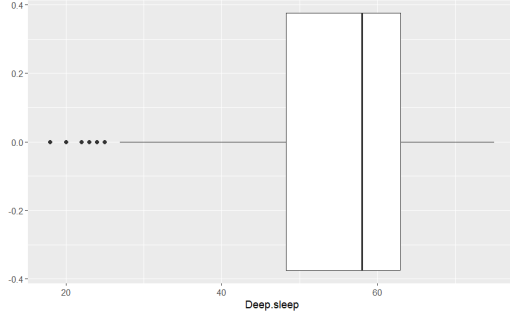## 2.2 Modeling Deep Sleep Percentage

### 2.2.1 Brief Introduction

In this part, we are to conduct a model to discover the relationships between deep sleeping percentage and other factors. Through model construction and analyses, we will find out what factors have more impacts. Furthermore, we can predict and conduct model evaluations based on the reduced model.
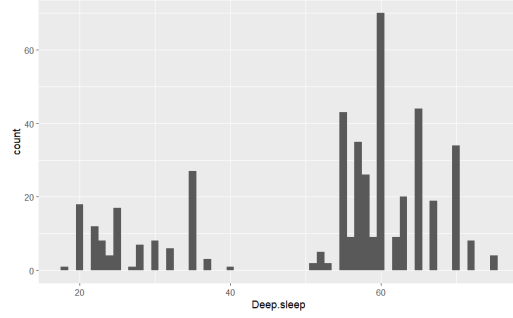
### 2.2.2 Data processing

Firstly, we conduct data exploration.

| mean | median | variance | standard deviation | min | max |
|---|---|---|---|---|---|
| 52.82301 | 58 | 245.0551 | 15.65424 | 18 | 75 |

Table 4. Statistics for deep sleep percentage
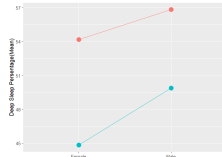


(a) Boxplot of Deep Sleep Percentage
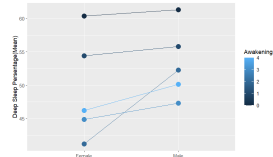


(b) Distribution

### 2.2.3 Base Model

From the figure, the distribution of deep sleep percentage is left-skewed. To figure out what factors contribute to the distribution, we are to conduct a base model.

```
## lm_ds1 = lm(Deep.sleep ~ Age + Gender + Awakenings + Caffeine.consumption +
## Alcohol.consumption + Smoking.status + Exercise.frequency, data = na.omit(data_ds))
## summary(lm_ds1)
```
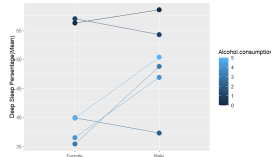
Based on the model, since parts of the factors are discrete and within a small range, we can conduct dummy variables as a method to observe the impact of some factors.
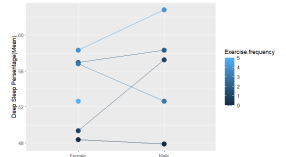


(c) Smoking



(d) Awakenings



(e) Alcohol



(f) Exercise

From the output, it seems that gender, smoking status, awakenings, alcohol consumption and exercise frequency all influence the percentage of deep sleep. As a result, we need to conduct a F-test to determine which reduced model is better in a precise way.

### 2.2.4 Reduced Model and F-test

We conduct a F-test to find the best reduced model.

```
## step(lm_ds1, test="F")
## Start:   AIC=2025.8
## Deep.sleep ~ Age + Gender + Awakenings + Caffeine.consumption +
## Alcohol.consumption + Smoking.status + Exercise.frequency
## ......
## Call:
## lm(formula = Deep.sleep ~ Gender + Awakenings + Alcohol.consumption +
##      Smoking.status + Exercise.frequency, data = na.omit(data_ds))
## Coefficients:
## (Intercept)  GenderMale  Awakenings  Alcohol.consumption  Smoking.statusYes
## Exercise.frequency
## 60.0406       2.8486      -2.8659     -3.0021              -6.5835            0.9648
```

Furthermore, we can discuss whether the reduced model is better using anova table.

```
## Model 1: Deep.sleep ~ Gender + Awakenings + Alcohol.consumption +
## Smoking.status + Exercise.frequency
## Model 2: Deep.sleep ~ Age + Gender + Awakenings + Caffeine.consumption +
## Alcohol.consumption + Smoking.status + Exercise.frequency
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1 382     69366
## 2 380     68934  2     432.05 1.1908 0.3051
```

From the table, it's safe for us to conclude that the reduced model is better. Based on the reduced model, we can conduct further analyses.

### 2.2.5 Statistical Inference and Predictions

Based on the reduced model, we can conduct predictions using R.

```
## confint(lm_ds2, level = 0.95)
## predict(lm_ds2, data.frame(Gender="Male", Awakenings=0,
## Alcohol.consumption=0, Smoking.status='No', Exercise.frequency=1), interval =
## "confidence")
## predict(lm_ds2, data.frame(Gender="Male", Awakenings=0,
## Alcohol.consumption=0, Smoking.status='No', Exercise.frequency=1), interval =
## "prediction")
```

```
                        2.5 %     97.5 %
(Intercept)         56.84120039 63.239924
GenderMale           0.02457003  5.672578
Awakenings          -3.92242949 -1.809390
Alcohol.consumption -3.86057417 -2.143622
Smoking.statusYes   -9.46502165 -3.702038
Exercise.frequency  -0.02672158  1.956261
```

```
        fit       lwr      upr
1 63.85391 60.55383 67.15398
        fit       lwr      upr
1 63.85391 37.15397 90.55384
```

(g) 95% Confidence interval      (h) Confidence intervals and Predictions

According to the 2 figures, we can predict that for a male who has a great lifestyle (no alcohol consumption, no smoking, exercise frequently) and doesn't wake up in the middle of the sleep, he is expected to enjoy about 63.85% of deep sleep every day, ranging from 60.55% to 67.15% at 95% confidence and predicted between 37.15% to 90.55%.

### 2.2.6 Testing for Collinearity

To check whether the regression assumptions are violated, we can check collinearity in R.

```{r}
vif(lm_ds2)
summary(lm_ds2)$r.squared
```

```
          Gender    Awakenings Alcohol.consumption  Smoking.status Exercise.frequency
        1.101965      1.131420            1.056793        1.033856           1.136019
[1] 0.2607818
```
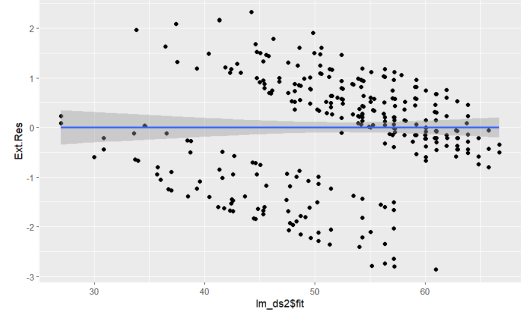
Figure 6. Code and Output of the test

According to the figure, the VIF values for all predictor variables are well below 5, indicating no significant multicollinearity in the model.

### 2.2.7 Model Evaluation

The data for model evaluation can be found when testing reduced models. With the reduced model, we can plot Residuals vs. Fitted Values.

Residual standard error: 13.48 on 382 degrees of freedom
Multiple R-squared:  0.2608,    Adjusted R-squared:  0.2511
F-statistic: 26.95 on 5 and 382 DF,  p-value: < 2.2e-16

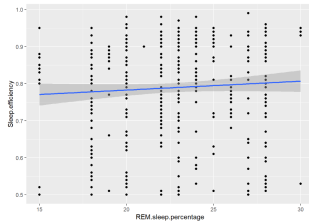(a) R output



(b) Residuals vs. Fitted Values

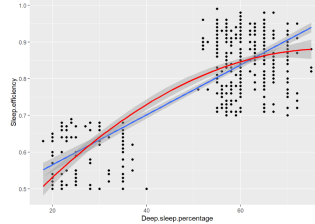## 2.3 Sleep Efficiency vs Sleep Stages

### 2.3.1 Intro

The conclusion for Sleep efficiency and Deep sleep percentage are quite similar. Therefore, we explore whether there are correlation between Sleep efficiency and percentages of the three sleep stages.
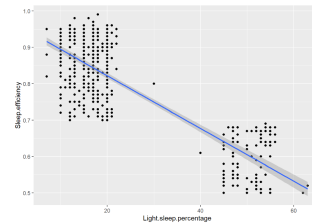
### 2.3.2 Data Exploration

The figures below show the scatter plot and fitted line (curve) of Sleep efficiency on individual stage percentage.



(c) REM sleep



(d) Deep sleep



(e) Light sleep

Figure 7. Efficiency vs individual stage percentage

According to figure 7, there may exists positive correlation between Sleep efficiency and Deep sleep percentage.

### 2.3.3 The Two-Predictor Model

The full model should be **Sleep.efficiency** $= \beta_0 + \beta_1$ **REM** $+ \beta_2$ **Deep** $+ \beta_3$ **Light** $+ \epsilon$

However, there may exists collinearity since the sum of the three stages' percentages should be 100. Therefore, we need to check collinearity between any two pairs.

According to the pairwise scatterplot and VIF (see table 5), we find strong negative correlation between Deep sleep percentage and Light sleep percentage. Therefore, we decided to drop Light sleep percentage to eliminate collinearity.

| Drop REM | | Drop light | |
|---|---|---|---|
| Deep | light | REM | Deep |
| 19.717 | 19.717 | 1.045293 | 1.045293 |

Table 5. vif result for two-predictor model

Then we get the reduced model: **Sleep efficiency** $= \beta_0 + \beta_1$ **REM** $+ \beta_2$ **Deep**

This model pass the F-test and two coefficients both pass t-test. And the boxcox plot shows we don't need to transform variables.
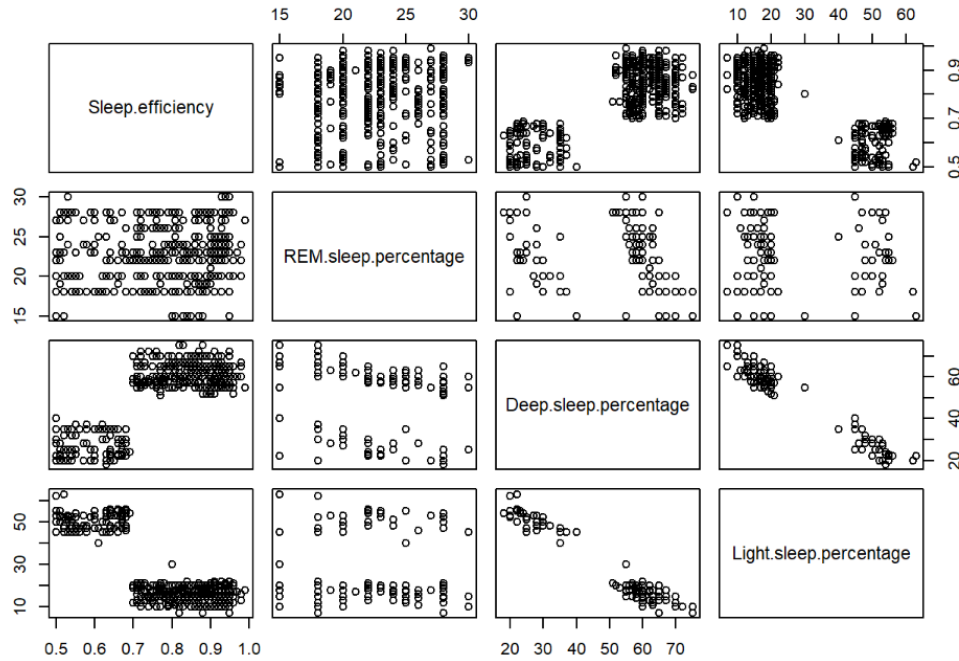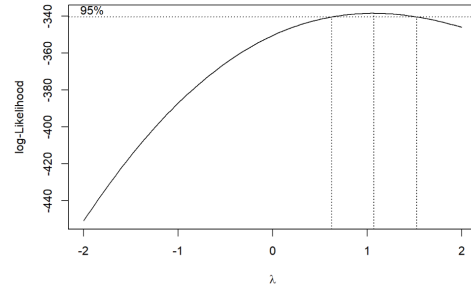
Figure 8. Pairwise scatterplots



```
## Coefficients:
##                      Estimate Std. Error t value Pr(>|t|)
## (Intercept)         0.2020206 0.0294779   6.853 2.39e-11 ***
## REM.sleep.percentage 0.0090709 0.0010576   8.577  < 2e-16 ***
## Deep.sleep.percentage 0.0072271 0.0002382 30.339  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07746 on 449 degrees of freedom
## Multiple R-squared:  0.6734,  Adjusted R-squared:  0.672
## F-statistic: 462.9 on 2 and 449 DF,  p-value: < 2.2e-16
```

(a) Summary  (b) Box Cox Plot

Figure 9. Summary and Box Cox Plot of the Two-Predictor Model

### 2.3.4 The reduced model

Since pairwise scatterplots don't show significant correlation between Sleep efficiency and REM sleep percentage, we wonder whether we can drop REM sleep percentage. The results of anova analysis are shown in table 6.

| Res.Df | RSS | Df | Sum of Sq | F | Pr( F) |
|--------|-----|-----|-----------|-----|--------|
| 450 | 3.135233 | | | | |
| 449 | 2.693872 | 1 | 0.4413616 | 73.56377 | 1.594602e-16 |

Table 6. Anova table for dropping REM

P-value is larger than 0.05, so we still should use the two-predictor model.

### 2.3.5 Regression Diagnostics

Residuals plot and QQ plot:



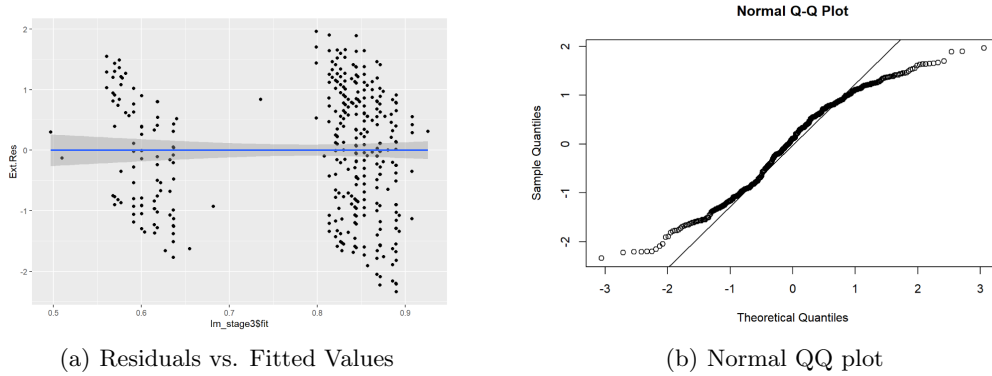(a) Residuals vs. Fitted Values                    (b) Normal QQ plot

Figure 10. Residual plot and QQ Plot

Residual plot doesn't show significant non-linearity, and residuals are randomly distributed. QQ plot is light tailed, so residuals are not normally distributed.

## 2.4 Result Summary

| Models | parameters | Adjusted R-sqaured | RMSE | F statistic | degrees of freedom |
|--------|-----------|--------------------|------|-------------|--------------------|
| Model1 | Table 2 | 0.5333 | 0.1604 | 89.46 | 382 |
| Model2 | R output 2.2.4 | 0.2511 | 13.48 | 26.95 | 382 |
| Model3 | Figure 9(a) | 0.672 | 0.07746 | 462.9 | 449 |

Table 7. Result summary

## 3. DISCUSSIONS

**Transformation of variables.** The transformation of variables is very important in data analysis, especially linear regression. In this report, it's almost impossible that the pattern in human's sleeping activity can be predicted by other variables linearly. Consequently, we can transform variables to try to build models more reliable.

**Model diagnostic.** We have an insight that the data is never perfect, which means that they may violate some of the assumptions of linear regression. So as we build our models step by step, we keep testing outliers, collinearity, normal distribution, etc.

**Graphic Method.** We also use many graphical methods in different stages when building models. Specifically, we find that the box plot and histogram can be used to show the distribution in data exploration, the pairwise scatter plot and QQ plot can be used to find whether the regression assumptions are violated. The residual plot can evaluate the model performance. Moreover, side-by-side figures combined with dummy variables are also very crucial to discover correlations.

**Feature selections.** We keep using automatic R feature selection like backward elimination, since there are countless combinations waiting for us to try and explore. The *step()* function in R can identify the most relevant predictors for our linear regression model. In addition, we always use **Anova table** to compare and contrast models.

## 4. CONCLUSIONS

In this report, we conducted a comprehensive data analysis and employed linear regression techniques to investigate the relationships in sleep. Our primary objectives are to identify patterns and correlations within a new sleeping dataset and to develop a predictive model using linear regression. We build three models: We start from building a relationship between sleep efficiency and several lifestyle variables. Then we focus on the deep sleep percentage, After that, we find these two models share some similarity in predictors, so we study the correlation between them. In general, this report utilizes linear regression to study the relationship between individuals' sleep efficiency, their health habits, as well as different stages of sleep including deep sleep.

## REFERENCES

1. Berry RB, Brooks R, Gamaldo CE, et al. (2012). The AASM manual for the scoring of sleep and associated events. Darien, IL: American Academy of Sleep Medicine.

2. Åkerstedt T, Hume K, Minors D, Waterhouse J. (1994). The meaning of good sleep: A longitudinal study of polysomnography and subjective sleep quality. J Sleep Res. 3:152–8.

3. Brand S, Gerber M, Kalak N, et al. Adolescents with greater mental toughness show higher sleep efficiency, more deep sleep and fewer awakenings after sleep onset[J]. Journal of Adolescent Health, 2014, 54(1): 109-113.

4. Hoffstein V, Mateika J, Mateika S. (1991). Snoring and sleep architecture. Am Rev Respir Dis. 143:92–6.

5. Janson C, De Backer W, Gislason T, et al. (1996). Increased prevalence of sleep disturbances and daytime sleepiness in subjects with bronchial asthma: A population study of young adults in three European countries. Eur Respir J. 9:2132–8.

6. Aydoğdu M, Ciftci B, Firat GS, et al. (2005). Assessment of sleep with polysomnography in patients with interstitial lung disease. Tuberk Toraks. 54:213–21.

7. Blackwell T, Yaffe K, Ancoli-Israel S, et al. (2006). Poor sleep is associated with impaired cognitive function in older women: The study of osteoporotic fractures. J Gerontol A Biol Sci Med. 61:405–10.

8. Scheer F, Zeitzer J, Ayas N, et al. (2006). Reduced sleep efficiency in cervical spinal cord injury; association with abolished night time melatonin secretion. Spinal Cord. 44:78–81.

9. Gruber R, Grizenko N, Schwartz G, et al. (2007). Performance on the continuous performance test in children with ADHD is associated with sleep efficiency. Sleep. 30:1003–9.

10. Yu J-M, Tseng I, Yuan R-Y, et al. (2009). Low sleep efficiency in patients with cognitive impairment. Acta Neurol Taiwan. 18:91–97.