

Homework 8

Jasmine Nakayama

April 12, 2018

Link to repository: (<https://github.com/jynakay/Assignments> (<https://github.com/jynakay/Assignments>))
[<https://github.com/jynakay/Assignments> (<https://github.com/jynakay/Assignments>)].

```
library(tidyverse)
library(class)
library(rpart)
library(NHANES)
library(RColorBrewer)
library(plot3D)
library(parallel)
library(randomForestSRC)
library(ggRandomForests)
library(mosaic)
library(dplyr)
```

Problem 1

Create the NHANES dataset again, just like we did in class, only using sleep trouble (variable name = SleepTrouble) as the dependent variable, instead of SleepTrouble.

What is the marginal distribution of sleep trouble?

```
#create dataset
people <- NHANES %>% dplyr::select(Age, Gender, SleepTrouble, BMI, HHIncome, PhysActive)
glimpse(people)
```

```
## Observations: 10,000
## Variables: 6
## $ Age          <int> 34, 34, 34, 4, 49, 9, 8, 45, 45, 45, 66, 58, 54, ...
## $ Gender       <fct> male, male, male, male, female, male, male, femal...
## $ SleepTrouble <fct> Yes, Yes, Yes, NA, Yes, NA, NA, No, No, No, No, N...
## $ BMI          <dbl> 32.22, 32.22, 32.22, 15.30, 30.57, 16.82, 20.64, ...
## $ HHIncome     <fct> 25000-34999, 25000-34999, 25000-34999, 20000-2499...
## $ PhysActive   <fct> No, No, No, NA, No, NA, NA, Yes, Yes, Yes, Yes, Y...
```

```
#marginal distribution of sleep trouble
tally(~ SleepTrouble, data = people, format = "percent")
```

```
## SleepTrouble
##    No    Yes <NA>
## 57.99 19.73 22.28
```

Recall from our prior work, the packages work better if the dataset is a dataframe, and the variables are numeric.

```
# Convert to dataframe
people <- as.data.frame(people)
class(people)
```

```
## [1] "data.frame"
```

```
glimpse(people)
```

```
## Observations: 10,000
## Variables: 6
## $ Age          <int> 34, 34, 34, 4, 49, 9, 8, 45, 45, 45, 66, 58, 54, ...
## $ Gender       <fct> male, male, male, male, female, male, male, femal...
## $ SleepTrouble <fct> Yes, Yes, Yes, NA, Yes, NA, NA, No, No, No, No, N...
## $ BMI          <dbl> 32.22, 32.22, 32.22, 15.30, 30.57, 16.82, 20.64, ...
## $ HHIncome     <fct> 25000-34999, 25000-34999, 25000-34999, 20000-2499...
## $ PhysActive   <fct> No, No, No, NA, No, NA, NA, Yes, Yes, Yes, Yes, Y...
```

```
#convert variables to numeric
people$Age <- as.numeric(people$Age)
people$Gender <- as.numeric(people$Gender)
people$SleepTrouble <- as.numeric(people$SleepTrouble)
people$BMI <- as.numeric(people$BMI)
people$HHIncome <- as.numeric(people$HHIncome)
people$PhysActive <- as.numeric(people$PhysActive)

people <- na.omit(people)

glimpse(people)
```

```
## Observations: 7,037
## Variables: 6
## $ Age          <dbl> 34, 34, 34, 49, 45, 45, 45, 66, 58, 54, 58, 50, 3...
## $ Gender       <dbl> 2, 2, 2, 1, 1, 1, 1, 2, 2, 2, 1, 2, 2, 2, 1, 1, 2...
## $ SleepTrouble <dbl> 2, 2, 2, 2, 1, 1, 1, 1, 1, 2, 1, 1, 1, 2, 1, 1, 2...
## $ BMI          <dbl> 32.22, 32.22, 32.22, 30.57, 27.24, 27.24, 27.24, ...
## $ HHIncome     <dbl> 6, 6, 6, 7, 11, 11, 11, 6, 12, 10, 11, 4, 6, 4, 1...
## $ PhysActive   <dbl> 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 1, 1, 2, 2, 2...
```

Problem 2

Apply the k-nearest neighbor procedure to predict SleepTrouble from the other covariates. Use $k = 1, 3, 5$, and 20 .

```
#Apply k-nearest neighbor approach to predict SleepTrouble for k = 1, 3, 5, 20
knn.1 <- knn(train = people, test = people, cl = people$SleepTrouble, k = 1)
knn.3 <- knn(train = people, test = people, cl = people$SleepTrouble, k = 3)
knn.5 <- knn(train = people, test = people, cl = people$SleepTrouble, k = 5)
knn.20 <- knn(train = people, test = people, cl = people$SleepTrouble, k = 20)
```

Problem 3

Now let's see how well these classifiers work overall.

```
# Calculate the percent predicted correctly

100*sum(people$SleepTrouble == knn.1)/length(knn.1)
```

```
## [1] 100
```

```
100*sum(people$SleepTrouble == knn.3)/length(knn.3)
```

```
## [1] 92.01364
```

```
100*sum(people$SleepTrouble == knn.5)/length(knn.5)
```

```
## [1] 88.68836
```

```
100*sum(people$SleepTrouble == knn.20)/length(knn.20)
```

```
## [1] 78.62726
```

Problem 4

What about success overall?

```
# Another way to look at success rate against increasing k

table(knn.1, people$SleepTrouble)
```

```
##
## knn.1    1    2
##      1 5239    0
##      2    0 1798
```

```
table(knn.3, people$SleepTrouble)
```

```
##
## knn.3    1    2
##      1 5062  385
##      2  177 1413
```

```
table(knn.5, people$SleepTrouble)
```

```
##  
## knn.5    1    2  
##      1 5029  586  
##      2  210 1212
```

```
table(knn.20, people$SleepTrouble)
```

```
##  
## knn.20    1    2  
##      1 5098 1363  
##      2  141  435
```