

# Assignment: Esoph and Young People Survey

*BOUN - ETM 58D*

*Due Date April 24, 2018*

This is an individual assignment. You are given 3 data sets to explore and build models. Add the assignment outcomes to your individual Progress Journals in 3 different html pages built with RMarkdown with all your code and outputs. **Please** do not forget to add detailed comments and explanations in proper content length to show that you understand these concepts and you can express your findings in a coherent way. Even though there is verbal guidance about the models you can use in the questions, you can use all the models you learned in the class (including random forests).

Collaboration is allowed but your work should be your own. Since copy pasting is so easy these days, high similarity between two submissions will result in penalized or nulled grades for this assignment. Those data sets are popular on internet. If you find an inspiration, please state it in a references section with links.

## Assignment 1: Esoph and Young People Survey Data

Questions below ask for your insight about the data. There is no single and constructed way to the right answer. Objective of this assignment is to help you help yourself to get data from outside resources, analyze the data, validate your model and convey your conclusions in a clear and reproducible way. You are expected to use R Markdown outputs on your Progress Journals to show your work.

1. Use `esoph` data set to see if (o)esophageal cancer is related to alcohol consumption, age or tobacco consumption. (Just write `esoph` to your R console to get the data.)
2. Use the Young People Survey's Hobbies & Interests category answers to infer "meaning" from the data. You are expected to use methods described in the lecture notes. But also you are welcome to use different methods as well. You can get the data from Kaggle.
  1. History: Not interested 1-2-3-4-5 Very interested (integer)
  2. Psychology: Not interested 1-2-3-4-5 Very interested (integer)
  3. Politics: Not interested 1-2-3-4-5 Very interested (integer)
  4. Mathematics: Not interested 1-2-3-4-5 Very interested (integer)
  5. Physics: Not interested 1-2-3-4-5 Very interested (integer)
  6. Internet: Not interested 1-2-3-4-5 Very interested (integer)
  7. PC Software, Hardware: Not interested 1-2-3-4-5 Very interested (integer)
  8. Economy, Management: Not interested 1-2-3-4-5 Very interested (integer)
  9. Biology: Not interested 1-2-3-4-5 Very interested (integer)
  10. Chemistry: Not interested 1-2-3-4-5 Very interested (integer)
  11. Poetry reading: Not interested 1-2-3-4-5 Very interested (integer)
  12. Geography: Not interested 1-2-3-4-5 Very interested (integer)
  13. Foreign languages: Not interested 1-2-3-4-5 Very interested (integer)
  14. Medicine: Not interested 1-2-3-4-5 Very interested (integer)
  15. Law: Not interested 1-2-3-4-5 Very interested (integer)
  16. Cars: Not interested 1-2-3-4-5 Very interested (integer)
  17. Art: Not interested 1-2-3-4-5 Very interested (integer)
  18. Religion: Not interested 1-2-3-4-5 Very interested (integer)
  19. Outdoor activities: Not interested 1-2-3-4-5 Very interested (integer)
  20. Dancing: Not interested 1-2-3-4-5 Very interested (integer)
  21. Playing musical instruments: Not interested 1-2-3-4-5 Very interested (integer)

22. Poetry writing: Not interested 1-2-3-4-5 Very interested (integer)
23. Sport and leisure activities: Not interested 1-2-3-4-5 Very interested (integer)
24. Sport at competitive level: Not interested 1-2-3-4-5 Very interested (integer)
25. Gardening: Not interested 1-2-3-4-5 Very interested (integer)
26. Celebrity lifestyle: Not interested 1-2-3-4-5 Very interested (integer)
27. Shopping: Not interested 1-2-3-4-5 Very interested (integer)
28. Science and technology: Not interested 1-2-3-4-5 Very interested (integer)
29. Theatre: Not interested 1-2-3-4-5 Very interested (integer)
30. Socializing: Not interested 1-2-3-4-5 Very interested (integer)
31. Adrenaline sports: Not interested 1-2-3-4-5 Very interested (integer)
32. Pets: Not interested 1-2-3-4-5 Very interested (integer)

## Assignment 2: Diamonds Data

Your assignment consists of finding the price of a diamond given its properties. You will use the `diamonds` data set in `ggplot2` package (which is inside `tidyverse`). You need to do your exploratory analysis well and come up with a predictive model. Your performance depends on the difference between the actual price of the diamond and the predicted price by the model. Use the `price` column as the response variable and other columns (except `diamond_id`) as predictors.

You are recommended to use CART but welcome to use any advanced method you like. Add your exploratory analysis to form a basis of your model and include references (with links) if you are inspired from similar analysis. Use the following code (and random seed) to form your train and test data. Remember, you should train your model on the train data and your real performance depends on the test data.

```
set.seed(503)
library(tidyverse)
diamonds_test <- diamonds %>% mutate(diamond_id = row_number()) %>%
  group_by(cut, color, clarity) %>% sample_frac(0.2) %>% ungroup()

diamonds_train <- anti_join(diamonds %>% mutate(diamond_id = row_number()),
  diamonds_test, by = "diamond_id")

diamonds_train
```

```
## # A tibble: 43,143 x 11
##   carat cut      color clarity depth table price     x     y     z
##   <dbl> <ord>   <ord> <ord>  <dbl> <dbl> <int> <dbl> <dbl> <dbl>
## 1 0.230 Ideal     E     SI2    61.5  55.0   326  3.95  3.98  2.43
## 2 0.210 Premium  E     SI1    59.8  61.0   326  3.89  3.84  2.31
## 3 0.230 Good     E     VS1    56.9  65.0   327  4.05  4.07  2.31
## 4 0.290 Premium  I     VS2    62.4  58.0   334  4.20  4.23  2.63
## 5 0.240 Very Good J     VVS2   62.8  57.0   336  3.94  3.96  2.48
## 6 0.240 Very Good I     VVS1   62.3  57.0   336  3.95  3.98  2.47
## 7 0.260 Very Good H     SI1    61.9  55.0   337  4.07  4.11  2.53
## 8 0.220 Fair     E     VS2    65.1  61.0   337  3.87  3.78  2.49
## 9 0.230 Very Good H     VS1    59.4  61.0   338  4.00  4.05  2.39
## 10 0.300 Good     J     SI1    64.0  55.0   339  4.25  4.28  2.73
## # ... with 43,133 more rows, and 1 more variable: diamond_id <int>
```

```
diamonds_test
```

```
## # A tibble: 10,797 x 11
##   carat cut    color clarity depth table price    x    y    z
##   <dbl> <ord> <ord> <ord>   <dbl> <dbl> <int> <dbl> <dbl> <dbl>
## 1 3.40 Fair D     I1      66.8  52.0 15964  9.42  9.34  6.27
## 2 0.900 Fair D     SI2     64.7  59.0  3205  6.09  5.99  3.91
## 3 0.950 Fair D     SI2     64.4  60.0  3384  6.06  6.02  3.89
## 4 1.00 Fair D     SI2     65.2  56.0  3634  6.27  6.21  4.07
## 5 0.700 Fair D     SI2     58.1  60.0  2358  5.79  5.82  3.37
## 6 1.04 Fair D     SI2     64.9  56.0  4398  6.39  6.34  4.13
## 7 0.700 Fair D     SI2     65.6  55.0  2167  5.59  5.50  3.64
## 8 1.03 Fair D     SI2     66.4  56.0  3743  6.31  6.19  4.15
## 9 1.10 Fair D     SI2     64.6  54.0  4725  6.56  6.49  4.22
## 10 2.01 Fair D     SI2     59.4  66.0 15627  8.20  8.17  4.86
## # ... with 10,787 more rows, and 1 more variable: diamond_id <int>
```

## Assignment 3: Spam Data

Original library is in UCI Database. See documentation on the website for further detail.

Your assignment consists of building a CART model to detect spam mail using UCI's Spambase data and analyze it. Your performance depends on correct specification of spam/non-spam mails in the test subset. You are going to use the RData file given on course webpage. Report your way of thinking, methodology, code and results.

You can load the data by using `load` command from your working directory or anywhere if you specify the path. For some installations, you can also double click the on the RData file to load. Name of the data frame is `spam_data` (same as the file name).

```
load("spam_data.RData")
head(spam_data)
```

Column names and short explanations are given below. For further details see the UCI documentation given in the above link.

train\_or\_test - 0 train, 1 test

spam\_or\_not - 0 not spam, 1 spam

V1 - word\_freq\_make

V2 - word\_freq\_address

V3 - word\_freq\_all

V4 - word\_freq\_3d

V5 - word\_freq\_our

V6 - word\_freq\_over

V7 - word\_freq\_remove

V8 - word\_freq\_internet

V9 - word\_freq\_order

V10 - word\_freq\_mail

V11 - word\_freq\_receive

V12 - word\_freq\_will  
V13 - word\_freq\_people  
V14 - word\_freq\_report  
V15 - word\_freq\_addresses  
V16 - word\_freq\_free  
V17 - word\_freq\_business  
V18 - word\_freq\_email  
V19 - word\_freq\_you  
V20 - word\_freq\_credit  
V21 - word\_freq\_your  
V22 - word\_freq\_font  
V23 - word\_freq\_000  
V24 - word\_freq\_money  
V25 - word\_freq\_hp  
V26 - word\_freq\_hpl  
V27 - word\_freq\_george  
V28 - word\_freq\_650  
V29 - word\_freq\_lab  
V30 - word\_freq\_labs  
V31 - word\_freq\_telnet  
V32 - word\_freq\_857  
V33 - word\_freq\_data  
V34 - word\_freq\_415  
V35 - word\_freq\_85  
V36 - word\_freq\_technology  
V37 - word\_freq\_1999  
V38 - word\_freq\_parts  
V39 - word\_freq\_pm  
V40 - word\_freq\_direct  
V41 - word\_freq\_cs  
V42 - word\_freq\_meeting  
V43 - word\_freq\_original  
V44 - word\_freq\_project  
V45 - word\_freq\_re  
V46 - word\_freq\_edu  
V47 - word\_freq\_table

V48 - word\_freq\_conference

V49 - char\_freq\_;

V50 - char\_freq\_(

V51 - char\_freq\_[

V52 - char\_freq\_!

V53 - char\_freq\_\$

V54 - char\_freq\_#

V55 - capital\_run\_length\_average

V56 - capital\_run\_length\_longest

V57 - capital\_run\_length\_total