

STA 314H1S: Problem Set 4

Except for question 4, the questions on this assignment are practice for the quiz on Friday, and are not to be handed in.

For question 4, please bring a printout of your output together with the source code. You will be asked to submit this printout at the end of the quiz.

1. (*Textbook Exercises*) Answer questions 1,2,4,5,6 from section 8.4, and questions 1,2,3 from section 9.7.
2. (*Support Vector Machine*) Recall the SVM classifier and the kernel trick that was used to construct it.
 - (a) Show that if $\mathbf{A} \in \mathbb{R}_{p \times p}$ is any PD matrix, then $K(\mathbf{z}_1, \mathbf{z}_2) = \mathbf{z}_1^T \mathbf{A} \mathbf{z}_2$ is a valid inner-product for $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^p$.
 - (b) Show that if $\mathbf{x} \in \mathbb{R}^p$ is a random vector with zero mean and covariance Σ_x , then $K(\mathbf{z}_1, \mathbf{z}_2) = \text{Cov}(\mathbf{z}_1^T \mathbf{x}, \mathbf{z}_2^T \mathbf{x})$ is a valid inner-product for $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^p$.
 - (c) For any $\mathbf{z} \in \mathbb{R}^p$, suppose $\Phi(\mathbf{z}) = (\phi_1(\mathbf{z}), \dots, \phi_q(\mathbf{z}))^T \in \mathbb{R}^q$, can we conclude that $K(\mathbf{z}_1, \mathbf{z}_2) = \Phi(\mathbf{z}_1)^T \Phi(\mathbf{z}_2)$ is a valid inner-product?
 - (d) Show that the kernel function $K(\mathbf{z}_1, \mathbf{z}_2) = (\mathbf{z}_1^T \mathbf{z}_2)^2$ that we learned in the lecture, is not a inner product.
 - (e) Describe why the Mercer's theorem implies the SVM is more useful in practice than the SVC with feature mapping?
 - (f) Describe the difference between the SVC and the maximal marginal classifier.
 - (g) Besides the budget parameter C , the choice of kernel in the SVM also controls the bias-variance tradeoff. Briefly explain why that is the case.
3. (*Decision Tree*) Suppose you collected data on the survival status of 6 passengers on the Titanic. The data includes their gender, age, weight, and height. The dataset is summarized in the table below.

Using this data, you built a decision tree to predict survival status.

Passenger	Survived	Gender	Age	Weight (kg)	Height (cm)
1	Yes	Male	22	70	170
2	No	Female	30	55	160
3	No	Male	45	78	175
4	Yes	Female	25	60	166
5	No	Male	55	80	172
6	No	Female	35	58	158

Table 1: Dataset of Titanic passengers

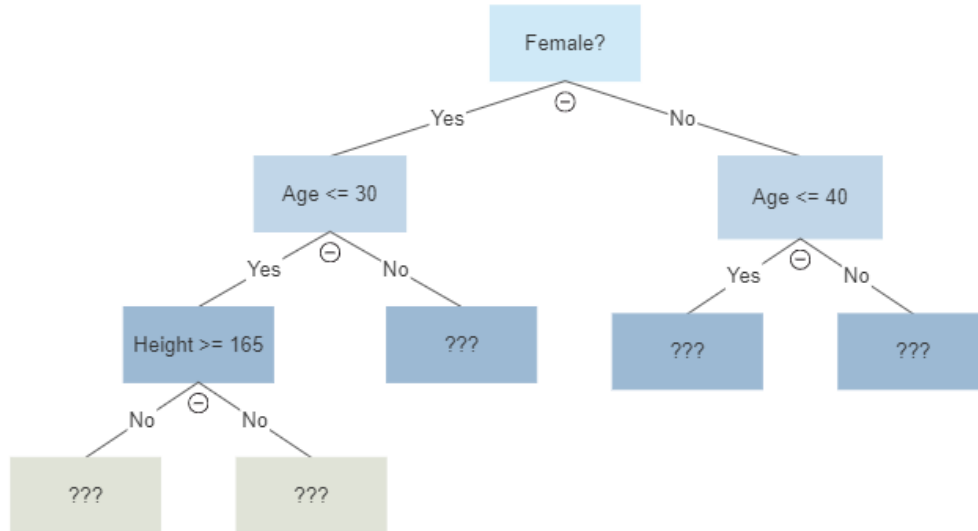


Figure 1: Decision tree for predicting survival on the Titanic.

- How many leaves does your decision tree have? How many observations are in each leaf?
 - Suppose your decision tree is built by minimizing the overall error rate in each leaf, fill in the missing classification results in each terminal leaf.
 - Predict the survival status of the following passengers based on your decision tree and explain your reasoning:
 - A 50-year-old male who weighs 76kg and is 172cm tall.
 - A 28-year-old female who weighs 58kg and is 160cm tall.
 - Critically discuss whether the decision tree might be prone to overfitting or underfitting based on the provided dataset and explain your reasoning.
4. (Computation) In this question, we will consider the `lung` data from the `survival` package in R. We are aiming to build a model to predict the survival time `time` using three features: `age`, `sex` (Male=1 Female=2) and `wt.loss` (the loss of weight in pounds, in the last six months). You have just received a new patient with feature $\mathbf{x}_{new} = (\text{age} = 21, \text{sex} = 2, \text{wt.loss} = 5)$.
- To start with, remove the observations with missing values, and separate the data into a training set with size 200, and the rest of the data to the testing set.
Some starter code to help you get started:

```

1 library(tidyverse)
2
3 ##### The dataset:
4 data <- survival::lung %>%
5   select(time, age, sex, wt.loss) %>%

```

```

6   drop_na()
7
8   ### Separation:
9   set.seed(123)
10  training_data <- #TODO
11  testing_data <- #TODO

```

- (b) Now, consider fitting a linear regression model with these features, and use it to predict the survival time for a female patient with `age = 21` and `wt.loss = 5`, called $\hat{y}(\mathbf{x}_{new})$.

Some starter code to help you get started:

```

1   ### Linear regression:
2   mod1 <- #TODO
3   new_data <- #TODO
4   mod1_pred <- predict(#TODO)
5   mod1_pred

```

- (c) To compute $SD(\hat{y}(\mathbf{x}_{new})) = \sqrt{\text{Var}(\hat{y}(\mathbf{x}_{new}))}$, we can no longer do it by simulating more datasets from the true data generating process as we have done in the PS1. However, the bootstrapping method allows us to estimate $SD(\hat{y}(\mathbf{x}_{new}))$ just using the single training set. Use $B = 100$ bootstrap samples to estimate $SD(\hat{y}(\mathbf{x}_{new}))$.

Some starter code to help you get started:

```

1   ### Bootstrapping for variance of 'mod1_pred'
2   B <- #TODO
3   ## The vector of yhat constructed from each bootstrap sample
4   yhat_vector <- c()
5   for (i in 1:B) {
6     index <- #TODO
7     training_data_i <- #TODO
8     boot_mod_i <- #TODO
9     yhat_vector[i] <- as.numeric(predict(boot_mod_i, new_data))
10  }
11  sd(yhat_vector)

```

- (d) Now, consider fitting a regression tree to this data.

```

1   ### Try regression tree:
2   library(tree)
3   tree_mod <- tree(time ~ age + sex + wt.loss, data = training_data
4     ,
5     split = "deviance",
6     control = tree.control(nobs = nrow(training_data), mincut = 1, minsize = 2, mindev = #TODO))
7   summary(tree_mod)
8   mean(summary(tree_mod)$residuals^2)

```

Try to vary the value of `mindev` from 1 to 0, and see how the structure of the tree changes, as well as how the residual sum of squares changes.

- (e) Now, suppose you want to use the regression tree above to obtain the prediction $\hat{y}(\mathbf{x}_{new})$. Estimate $SD(\hat{y}(\mathbf{x}_{new}))$ of your regression tree using $B = 100$ bootstrap samples. Consider both $mindev = 1$ and $mindev = 0$, what are their estimated $SD(\hat{y}(\mathbf{x}_{new}))$ respectively?

Some starter code to help you get started:

```

1 set.seed(123)
2 ### Bootstrapping for variance of a tree with mindev = 0
3 B <- 100
4 ## The vector of yhat constructed from each bootstrap sample
5 yhat_vector1 <- c()
6 for (i in 1:B) {
7   training_data_i <- #TODO
8   tree_mod_i <- tree(time ~ age + sex + wt.loss, data = #TODO,
9                     split = "deviance",
10                    control = tree.control(nobs = nrow(training_
11      data_i), mincut = 1, minsize = 2, mindev = 0))
12   yhat_vector1[i] <- as.numeric(predict(tree_mod_i, new_data))
13 }
14 sd(yhat_vector1)
15
16 ### Bootstrapping for variance of a tree with mindev = 1
17 B <- 100
18 ## The vector of yhat constructed from each bootstrap sample
19 yhat_vector2 <- c()
20 for (i in 1:B) {
21   training_data_i <- #TODO
22   tree_mod_i <- tree(time ~ age + sex + wt.loss, data = #TODO,
23                     split = "deviance",
24                    control = tree.control(nobs = nrow(training_
25      data_i), mincut = 1, minsize = 2, mindev = 1))
26   yhat_vector2[i] <- as.numeric(predict(tree_mod_i, new_data))
27 }
28 sd(yhat_vector2)

```

- (f) Suppose you want to use Bagging to construct an ensemble prediction using multiple trees, complete the following R function:

```

1 ### Bagging:
2 bagging <- function(tree_num = 100, mindev = 0, training_data,
3   testing_data){
4   # Initialize the prediction result yhat_vector
5   yhat_vector <- matrix(nrow = nrow(testing_data), ncol = tree_
6     num)
7   for (i in 1:tree_num) {
8     # Compute the prediction of each bootstrapped tree
9     training_data_i <- # TODO
10    tree_mod_i <- tree(time ~ age + sex + wt.loss, data =
11      training_data_i,
12                      split = "deviance",
13                      control = tree.control(nobs = nrow(
14        training_data_i), mincut = 1, minsize = 2, mindev = mindev))

```

```

11     yhat_vector[,i] <- as.numeric(predict(tree_mod_i, testing_
12     data))
13 }
14 ### Compute the final estimate as the sample mean
15 return(apply(#TODO))

```

- (g) Using your R function above, estimate $SD(\hat{y}(\mathbf{x}_{new}))$ of your Bagging prediction where `tree_num` = 100 and `mindev` = 0, using $B = 100$ bootstrap samples. Compare your answer to part (e).
- (h) Compute the empirical testing MSE of your bagging prediction from (g) with `tree_num` = 100 and `mindev` = 0, of the single regression tree with `mindev` = 0 from part (e), and of the linear regression model from part (b), using the reserved testing data.