

To begin with some preprocessing.

```
1 library(rpart)
2 library(rpart.plot)
3 library(magrittr)
4 library(dplyr)
5
6 dataset <- read.csv("insurance.csv")
7 dataset <- na.omit(dataset)
8 # Split the dataset
9 train_indices<- sample(1:nrow(dataset), size = 0.8*nrow(dataset))
10 train_set <- dataset[train_indices,]
11 test_set <- dataset[-train_indices,]
```

Question 1. (a) Create an OLS regression model.

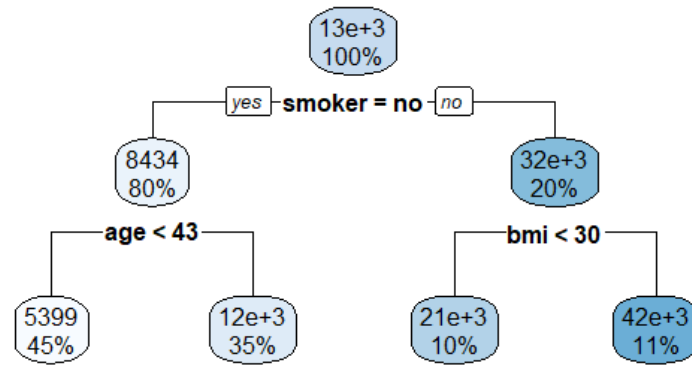
```
1 # Question 1 (a)
2 ins_lm <- lm(charges ~ age + factor(sex) + bmi + factor(smoker) +
3             factor(region), data=dataset)
4 report_lm <- summary(ins_lm)
5 write.table(round(report_lm$coefficients, digits=5),
6             file="1a.csv",
7             sep = ",",
8             col.names=NA)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-11556.95734	985.62619	-11.7255	0
age	258.5397	11.93749	21.65779	0
factor(sex)male	-111.57001	334.25839	-0.33378	0.7386
bmi	340.45918	28.71413	11.85685	0
factor(smoker)yes	23862.90634	414.82279	57.52554	0
factor(region)northwest	-304.10359	478.01318	-0.63618	0.52477
factor(region)southeast	-1039.2021	480.64936	-2.16208	0.03079
factor(region)southwest	-916.44143	479.72079	-1.91036	0.0563

(b) Create a decision (regression) tree.

```
1 library(rpart)
2 library(rpart.plot)
3 library(magrittr)
4 library(dplyr)
5
6 dataset <- read.csv("insurance.csv")
7 dataset <- na.omit(dataset)
8 # Split the dataset
9 train_indices<- sample(1:nrow(dataset), size = 0.8*nrow(dataset))
10 train_set <- dataset[train_indices,]
11 test_set <- dataset[-train_indices,]
```

(i) Plot a visual representation of the tree



- (ii) The depth is 2.
- (iii) There are 4 leaf groups
- (iv) As the following table shows.

smoker, age<43	smoker, age>=43	non-smoker, bmi<30	non-smoker, bmi>=30
28675	36716	7977	8843

- (v) Please refer to the plot in (i).

Question 2. The codes to compute the RMSE:

```

1 # Question 2
2 mse_oos <- function(actuals, preds) { # MSEs in this homework are all RMSE
3   sqrt(mean( (actuals - preds)^2 ))
4 }
5
6 mse_oos(dataset$charges, predict(ins_lm, dataset))
7 mse_oos(dataset$charges, predict(ins_tree, dataset))

```

- (a) The $RMSE_{oos}$ for the OLS regression model is 6068.683.
- (b) The $RMSE_{oos}$ for the decision tree model is 4321.024.

Question 3. (a) Thanks to 108071021 for help of this part (3(a)) of codes

```

1 # Question 3 (a)
2 # Thanks to 108071021 for help of this part of codes
3 bagged_retrain <- function(model, dataset, b){
4   resample <- unique(sample(1:nrow(dataset), replace = TRUE))
5   train_data <- dataset[resample,]
6   train_model <- update(model, data = train_data)
7   train_model
8 }
9
10 bagged_learn <- function(model, dataset, b=100){
11   lapply(1:b, bagged_retrain, model = model, dataset = dataset)
12 }
13
14 pred <- function(model, dataset, b){
15   model = model[[b]]
16   predict(model, dataset)
17 }
18
19 bagged_predict <- function(bagged_model, dataset, b){

```

```

20 prediction <- lapply(1:b, pred, model = bagged_model, dataset = dataset)
21 mse_oos(unlist(prediction), rep(unlist(dataset[7]), times = b))
22 }
23
24 # Question 3 (b)
25 lm_bagged_models <- bagged_learn(ins_lm, train_set, 100)
26 lm_bagged_mse <- bagged_predict(lm_bagged_models, test_set, 100)
27
28 # Question 3 (c)
29 ins_tree_models <- bagged_learn(ins_tree, train_set, 100)
30 ins_tree_bagged_mse <- bagged_predict(ins_tree_models, test_set, 100)

```

(b) The $RMSE_{oos}$ for the bagged OLS regression model is 6654.777.

(c) The $RMSE_{oos}$ for the bagged decision tree model is 5073.554. ■

Question 4. Thanks to 傅奕軒's code in the Microsoft Teams.

```

1 # Question 4 (a)
2 boosted_learn <- function(model, dataset, outcome, n=100, rate=0.1, type) {
3   # get data frame of only predictor variables
4   predictors <- dataset[, -which(names(dataset) %in% outcome)]
5   # Initialize residuals and models
6   res <- dataset[,outcome] # get vector of actuals to start
7   models <- list()
8   for (i in 1:n) {
9     this_model <- update(model, data = cbind(charges = res, predictors))
10    # update residuals with learning rate
11    if (type == "1")
12    {
13      res <- res - rate * this_model$fitted.values
14    }
15    else
16    {
17      res <- res - rate * this_model$y
18    }
19    models[[i]] <- this_model # Store model
20  }
21  list(models=models, rate=rate)
22 }
23
24 boost_predict <- function(boosted_learning, new_data) {
25   boosted_models <- boosted_learning$models
26   rate <- boosted_learning$rate
27   n <- length(boosted_models)
28
29   # get predictions of new_data from each model
30   predictions = lapply(1:n, function(i){
31     rate*predict(boosted_models[[i]], new_data)
32   })
33
34   pred_frame = as.data.frame(predictions) #>% unname
35   pred <- apply(pred_frame, 1, sum)
36   mse_oos(pred, new_data[,7]) # 7 for charges
37 }
38
39 # Question 4 (b)

```

```

40 ins_lm_boosted_model <- boosted_learn(ins_lm, train_set, outcome = "charges",
41                                     type = "l")
42 ins_lm_boosted_mse <- boost_predict(ins_lm_boosted_model, test_set)
43
44 # Question 4 (c)
45 ins_tree_stump <- rpart(charges ~ age + sex + bmi + smoker + region,
46                        data=dataset, cp=0, maxdepth=1)
47 ins_tree_boosted_model <- boosted_learn(ins_tree_stump, train_set,
48                                       outcome="charges", type='t')
49 ins_tree_boost_mse <- boost_predict(ins_tree_boosted_model, test_set)

```

(b) The $RMSE_{\text{oos}}$ for the boosted OLS regression model is 6638.21.

(c) The $RMSE_{\text{oos}}$ for the boosted decision tree model is 5166.023.

Note that the RMSE boosted method varies a lot. ■

Question 5. (a) The bagging of the decision tree

```

1  # Question 5 (a)
2  flag <- 0 # break the loop or not
3  maxdepth_ins <- 1 # maxdepth of the tree
4  maxdepth_vector <- c()
5  bagged_mse <- c()
6
7  while(flag == 0){
8    control <- rpart.control(maxdepth = maxdepth_ins, cp = 0)
9
10   # Using code in Question 3 (c)
11   ins_tree <- rpart(charges ~ bmi + age + sex + children + smoker + region,
12                   data = train_set, control = control)
13   ins_tree_models <- bagged_learn(ins_tree, train_set, 100)
14   ins_tree_bagged_mse <- bagged_predict(ins_tree_models, test_set, 100)
15
16   # append the result to the vector
17   maxdepth_vector <- c(maxdepth_ins)
18   bagged_mse <- c(bagged_mse, ins_tree_bagged_mse)
19
20   # To determine to brake the loop or not
21   if(length(bagged_mse) >= 2){
22     if(bagged_mse[maxdepth_ins-1] < bagged_mse[maxdepth_ins]){
23       flag <- 1 # break
24     }
25   }
26   maxdepth_ins <- maxdepth_ins + 1
27 }
28
29 result_dataframe_5a <- cbind(maxdepth_vector, bagged_mse)
30 names(result_dataframe_5a) <- c("maxdepth", "bagged_mse")
31 write.table(result_dataframe_5a, file="5a.csv", sep = ",", col.names=NA)

```

maxdepth_vector	bagged_mse
1	7108.743
2	4775.779
3	4572.248
4	4473.494
5	4512.021

(b) The boosting of the decision tree

```

1  # Question 5 (b)
2  flag <- 0 # break the loop or not
3  maxdepth_ins <- 1 # maxdepth of the tree
4  maxdepth_vector <- c()
5  boost_mse <- c()
6
7  while(flag == 0){
8      control <- rpart.control(maxdepth = maxdepth_ins, cp = 0)
9
10     # Using code in Question 4 (c)
11     ins_tree_stump <- rpart(charges ~ age + sex + bmi + smoker + region,
12                             data=dataset, control = control)
13     ins_tree_boosted_model <- boosted_learn(ins_tree_stump, train_set,
14                                             outcome="charges", type='t')
15     ins_tree_boost_mse <- boost_predict(ins_tree_boosted_model, test_set)
16
17     # append the result to the vector
18     maxdepth_vector <- c(maxdepth_ins)
19     boost_mse <- c(boost_mse, ins_tree_boost_mse)
20
21     # To determine to brake the loop or not
22     if(length(boost_mse) >= 2){
23         if(boost_mse[maxdepth_ins-1] < boost_mse[maxdepth_ins]){
24             flag <- 1 # break
25         }
26     }
27     maxdepth_ins <- maxdepth_ins + 1
28 }
29
30 result_dataframe_5b <- cbind(maxdepth_vector, boost_mse)
31 names(result_dataframe_5b) <- c("maxdepth", "bagged_mse")
32 write.table(result_dataframe_5b, file="5b.csv", sep = ",", col.names=NA)

```

maxdepth_vector	boost_mse
1	7095.81
2	4730.264
3	4481.213
4	4414.42
5	4471.158

