5. *What was your chosen prediction machine?* **Paste the code that produced your predicted values, including all values of tuning parameters if any, random number seeds, and explaining any variable names that are not obvious.** I should be able to run your code and produce the same results (or extremely similar if randomization is used). If I try and it doesn't work, there will be a major deduction.

➔ I choose randomforest as my prediction model.

```
 5  library(dplyr)
 6  library(MASS)     # For ridge regression
 7  library(glmnet)   # For LASSO
 8  library(nnet)     # Fits neural net models
 9  library(rgl)
10  library(FNN)
11  library(car)      # For ANOVA after logistic regression with nnet
12  library(mgcv)     # For GAM
13  library(klaR)     # For naive Bayes
14  library(rpart)          # For fitting classification trees
15  library(rpart.plot)     # For plotting classification trees
16  library(randomForest)   # For random forests
17  library(gbm)            # For boosting
18  source("Helper Functions.R")
19
20
21  ### Some of our code will be random, so we have to set the seed.
22  ### Use Mersenne-Twister for compatibility.
23  set.seed(46685326, kind = "Mersenne-Twister")
24  data = read.csv("P2Data2020.csv")
```

```r
317  ##############################################  Random Forest
318  data.rf = data
319  data.rf$Y = factor(data.rf$Y)
320
321  data.train.rf = data.rf[ind.random <= n.train,]
322  data.valid.rf = data.rf[ind.random > n.train,]
323  Y.train.rf = data.train.rf$Y
324  Y.valid.rf = data.valid.rf$Y
325
326  ### Set tuning parameters
327  all.mtrys = 1:6
328  all.nodesizes = c(1, 5, 10, 15, 20)
329  all.pars.rf = expand.grid(mtry = all.mtrys, nodesize = all.nodesizes)
330  n.pars = nrow(all.pars.rf)
331
332  M = 5 # Number of times to repeat RF fitting. I.e. Number of OOB errors
333
334  all.OOB.rf = array(0, dim = c(M, n.pars))
335  names.pars = apply(all.pars.rf, 1, paste0, collapse = "-")
336  colnames(all.OOB.rf) = names.pars
337  |
338  for(i in 1:n.pars){
339    ### Progress update
340    print(paste0(i, " of ", n.pars))
341
342    ### Get tuning parameters for this iteration
343    this.mtry = all.pars.rf[i, "mtry"]
344    this.nodesize = all.pars.rf[i, "nodesize"]
345
346    for(j in 1:M){
347      ### Fit RF, then get and store OOB errors
348      this.fit.rf = randomForest(Y ~ ., data = data.train.rf,
349                                 mtry = this.mtry, nodesize = this.nodesize)
350
351      pred.this.rf = predict(this.fit.rf)
352      this.err.rf = mean(Y.train.rf != pred.this.rf)
353
354      all.OOB.rf[j, i] = this.err.rf
355    }
356  }
357
358  ### Make a regular and relative boxplot
359  boxplot(all.OOB.rf, las=2, main = "OOB Boxplot")
360  rel.OOB.rf = apply(all.OOB.rf, 1, function(W) W/min(W))
361  boxplot(t(rel.OOB.rf), las=2,  # las sets the axis label orientation
362          main = "Relative OOB Boxplot") #3-10 has the lowest value
363
364  ### Best model looks like mtry = 3 and nodesize = 10.
365  fit.rf = randomForest(Y ~ ., data = data.train.rf,
366                        mtry = 3, nodesize = 10)
367
368  ### Get predictions and evaluate performance
369  pred.rf = predict(fit.rf, data.valid.rf)
370
371  table(Y.valid, pred.rf, dnn = c("Obs", "Pred"))
372  (mis.rf = mean(Y.valid != pred.rf))    #0.204
```

```r
374 ############################################## Data Prediction
375 data2 = read.csv("P2Data2020testX.csv")
376 fit.rf = randomForest(Y ~ ., data = data.train.rf,
377                       mtry = 3, nodesize = 10)
378
379 pred.rf = predict(fit.rf, data2)
380 write.table(pred.rf, "prediction1.csv", sep = ",", row.names = F, col.names = F)
```