

Homework 4–Tree based methods

Ngoc Duong

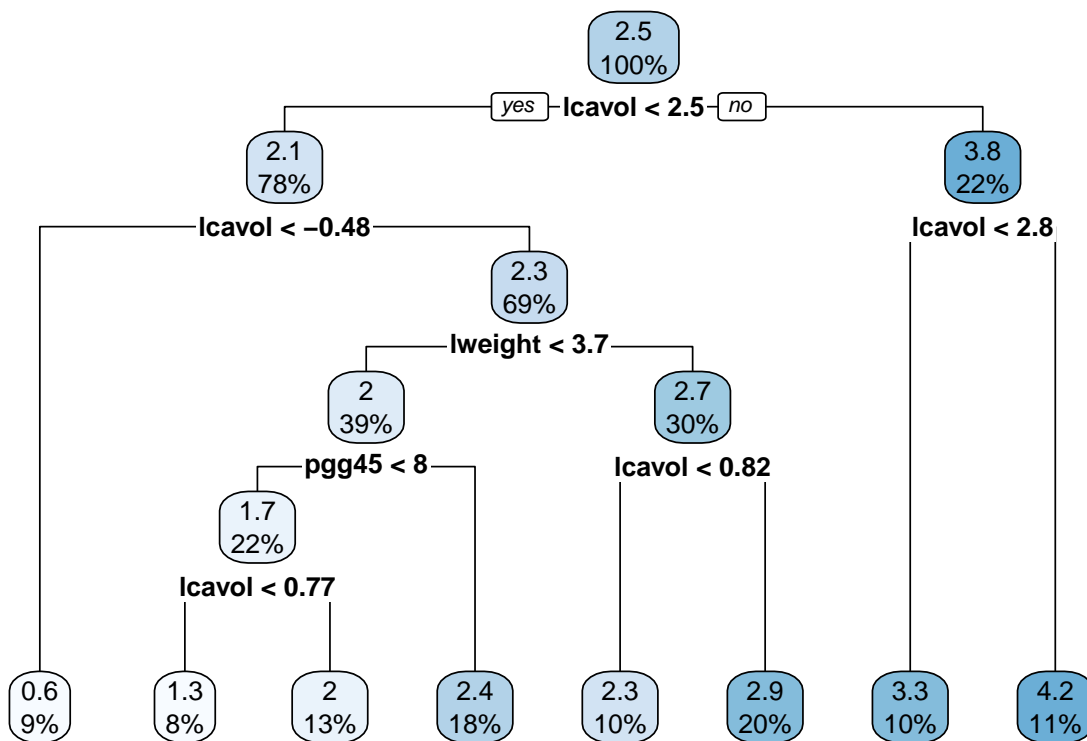
4/22/2020

Problem 1

```
data("Prostate") #import data
```

a. Fit a regression tree with “lpsa” as the response and the other variables as predictors

```
set.seed(13)
tree1 <- rpart(formula = lpsa ~ ., data = Prostate)
rpart.plot(tree1)
```



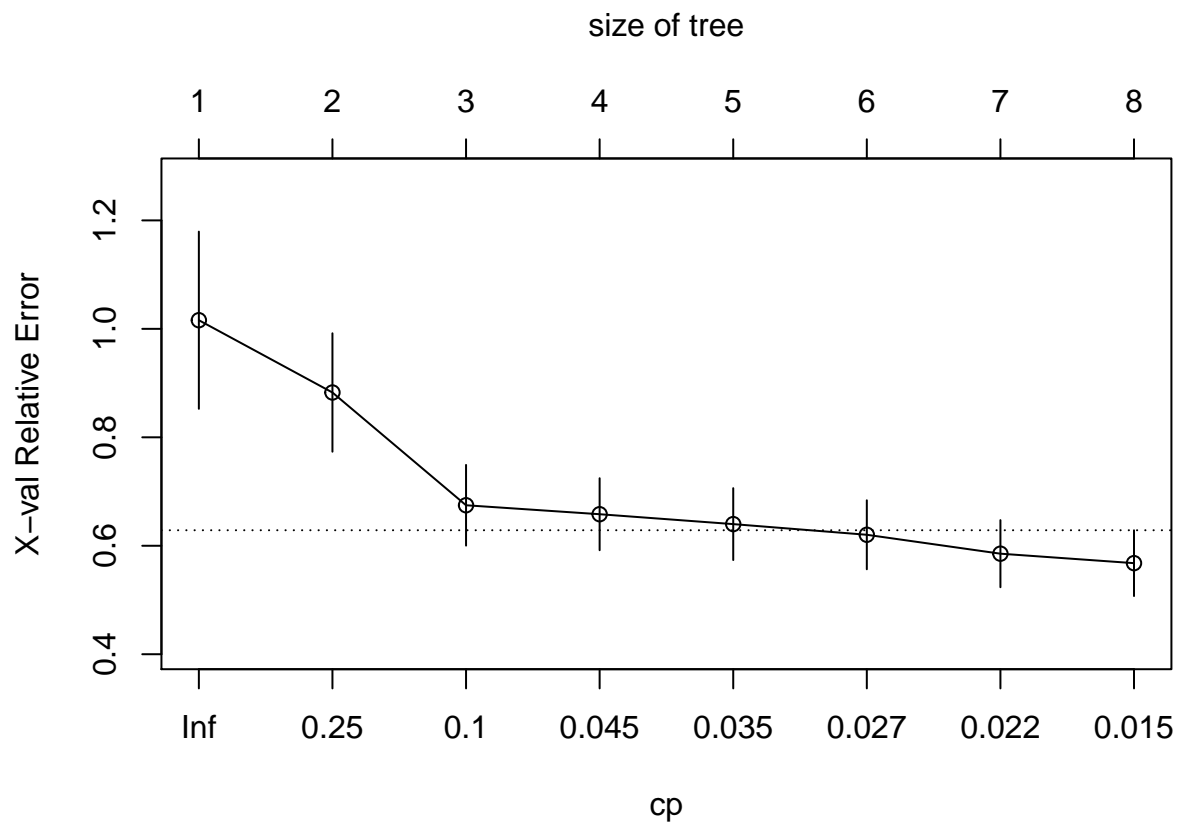
Next, we can use cross-validation to determine the optimal tree.

```
cpTable <- printcp(tree1)
```

```
##
## Regression tree:
## rpart(formula = lpsa ~ ., data = Prostate)
##
## Variables actually used in tree construction:
## [1] lcavol lweight pgg45
```

```
##
## Root node error: 127.92/97 = 1.3187
##
## n= 97
##
##      CP nsplit rel error  xerror   xstd
## 1 0.347108      0   1.00000 1.01591 0.163388
## 2 0.184647      1   0.65289 0.88262 0.109196
## 3 0.059316      2   0.46824 0.67467 0.074444
## 4 0.034756      3   0.40893 0.65816 0.066401
## 5 0.034609      4   0.37417 0.63994 0.066280
## 6 0.021564      5   0.33956 0.62026 0.063629
## 7 0.021470      6   0.31800 0.58540 0.061903
## 8 0.010000      7   0.29653 0.56793 0.060693
```

```
plotcp(tree1)
```



We can prune the tree based on the cp table

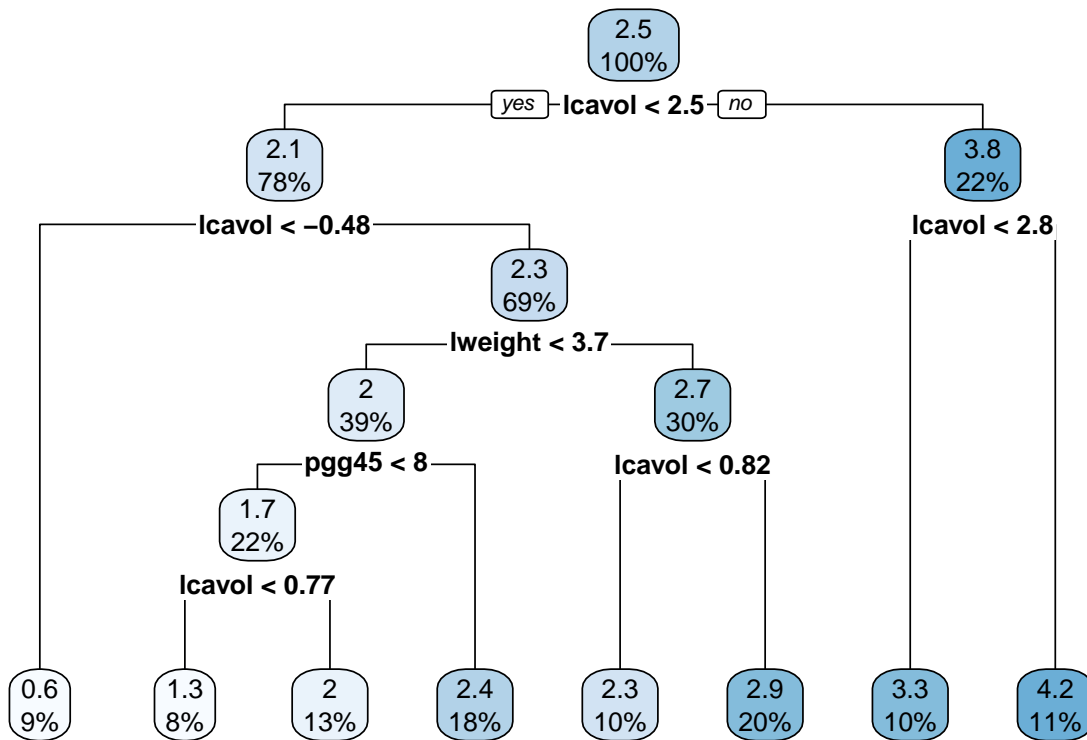
```
minErr <- which.min(cpTable[,4]) # minimum cross-validation error
tree2 <- prune(tree1, cp = cpTable[minErr,1])

tree3 <- prune(tree1, cp = cpTable[cpTable[,4] < cpTable[minErr,4] + cpTable[minErr,5],1][1]) # 1SE rule
```

The tree size corresponding to the lowest cross-validation error is 7 (splits), compared to 5 splits in the tree obtained by the 1SE rule.

b. Plot the final tree choice (7 splits), and interpret a terminal node

```
rpart.plot(tree2)
```



Looking at the fifth terminal node from the left, we can interpret it as: the predicted log(prostate specific antigen) is 2.3 for log(cancer volume) between -0.48 and 0.82 and log(prostate weight) more than 3.7. The number of observations falling into this branch and have predicted value 2.3 is 10% of the total observations.

c. Perform bagging and report the variable importance

```
set.seed(13)
bagging = ranger(lpsa ~ ., data = Prostate, mtry = 8,
  splitrule = "variance", min.node.size = 5,
  importance = "permutation",
  scale.permutation.importance = TRUE)
```

We can look at the lpsa prediction for first 10 subjects in the dataset using model obtained from bagging

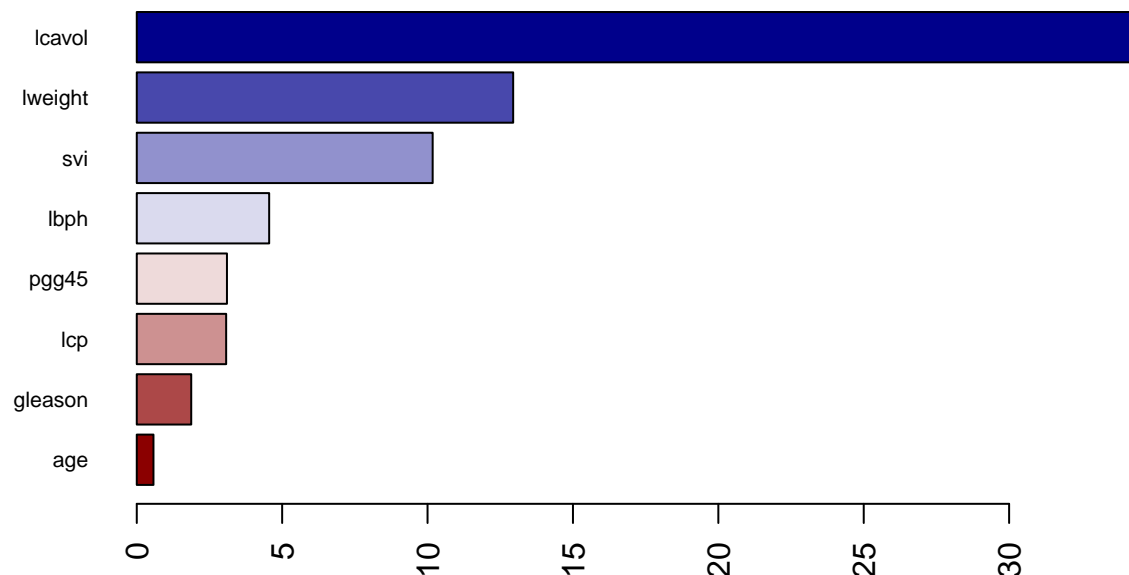
```
predict(bagging, data = Prostate[1:10,])$predictions
```

```
## [1] 0.3583673 0.1639139 0.6539709 0.2769208 0.8896994 0.3951543 0.9698987
## [8] 1.1680926 0.9338197 1.3424346
```

```
#variable importance
```

```
bagging2 = ranger(lpsa ~ ., data = Prostate, mtry = 8,
  splitrule = "variance", min.node.size = 5,
  importance = "permutation",
  scale.permutation.importance = TRUE)
```

```
barplot(sort(ranger::importance(bagging2), decreasing = FALSE),
  las = 2, horiz = TRUE, cex.names = 0.7,
  col = colorRampPalette(colors = c("darkred", "white", "darkblue"))(8))
```



The variable importance above was computed from permuting out-of-bag data, suggesting log(cancer volume) and log(prostate weight) have the most influence in the fitted models. This can also be seen from the chosen decision tree in part b. On the other hand, age, log(capsular penetration), and log(benign prostatic hyperplasia amount) carry the least importance in the fitted models.

d. Perform random forest and report the variable importance

```
set.seed(13)
#fast implementation
rf = ranger(lpsa ~ ., data = Prostate, mtry = 3, #decorrelate by picking mtry = 1/3 mtry in bagging
            splitrule = "variance", min.node.size = 5,
            importance = "permutation",
            scale.permutation.importance = TRUE)
```

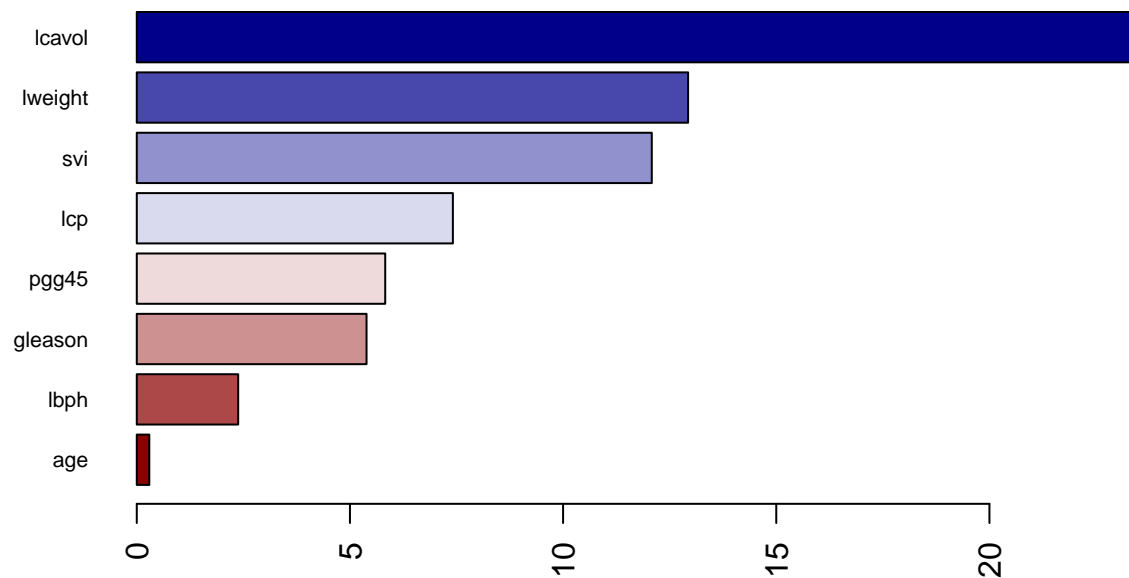
We can also look at the lpsa prediction for first 10 subjects in the dataset using model obtained from random forest

```
#prediction
predict(rf, data = Prostate[1:10,])$predictions
```

```
## [1] 0.4432477 0.3049250 0.8382494 0.3435001 0.9810472 0.5051814 1.1225813
## [8] 1.3338630 0.9114729 1.3479017
```

```
#importance importance
rf2 = ranger(lpsa ~ ., data = Prostate, mtry = 3, #decorrelate by picking mtry = 1/3 mtry in bagging
            splitrule = "variance", min.node.size = 5,
            importance = "permutation",
            scale.permutation.importance = TRUE)

barplot(sort(ranger::importance(rf2), decreasing = FALSE),
        las = 2, horiz = TRUE, cex.names = 0.7,
        col = colorRampPalette(colors = c("darkred", "white", "darkblue"))(8))
```



From the variable importance plot, we can see $\log(\text{cancer volume})$ and $\log(\text{prostate weight})$ have the most influence in the fitted models. The plot also suggests age, $\log(\text{capsular penetration})$, and $\log(\text{benign prostatic hyperplasia amount})$ carry the least importance in the fitted models. This is similar to the variable importance plot obtained from bagging models above.

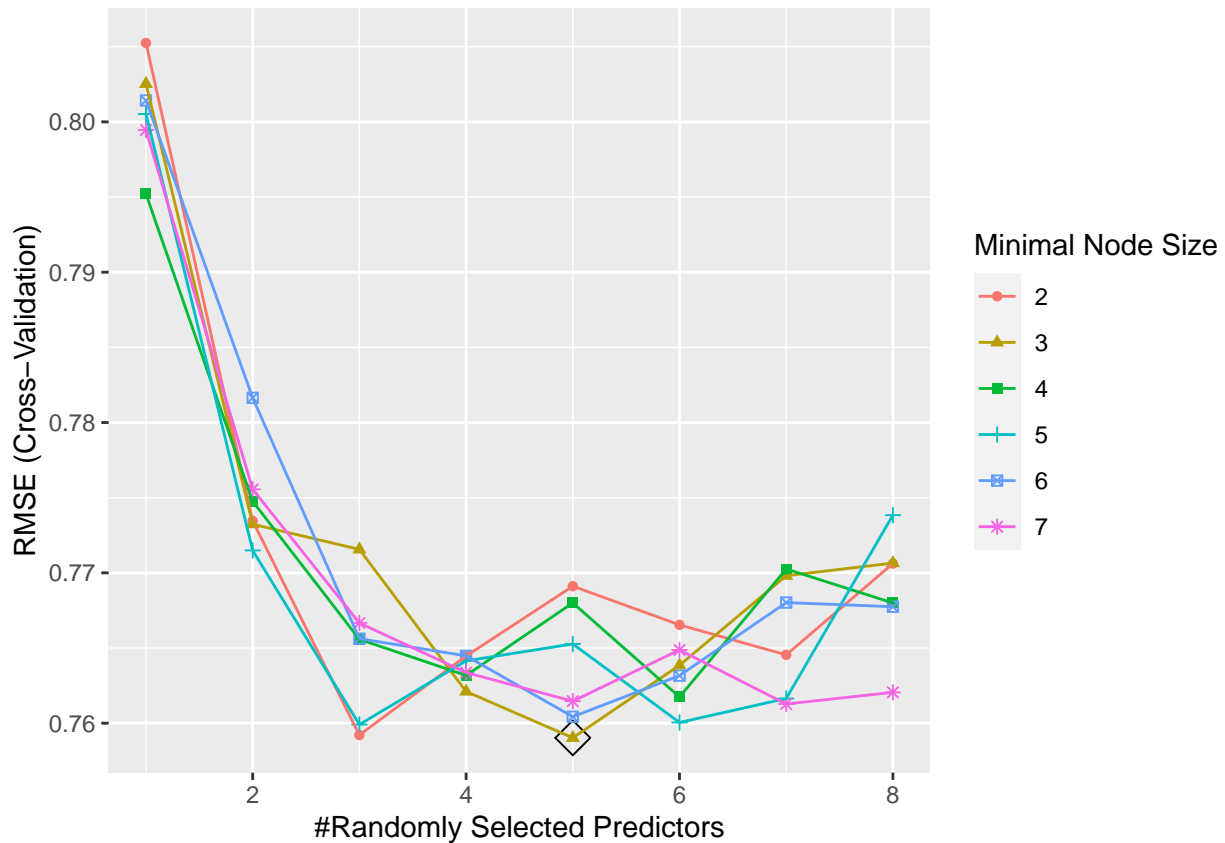
We can also try doing a grid search by caret for Random Forest

```
ctrl = trainControl(method = "cv")

rf.grid = expand.grid(mtry = 1:8,
                     splitrule = "variance",
                     min.node.size = 2:7)

set.seed(13)
rf.fit <- train(lpsa~., Prostate,
               method = "ranger",
               tuneGrid = rf.grid,
               trControl = ctrl)

ggplot(rf.fit, highlight = TRUE)
```



The tuned RF model has the best minimal node size of 3 and the randomly selected predictors mtry of 4.

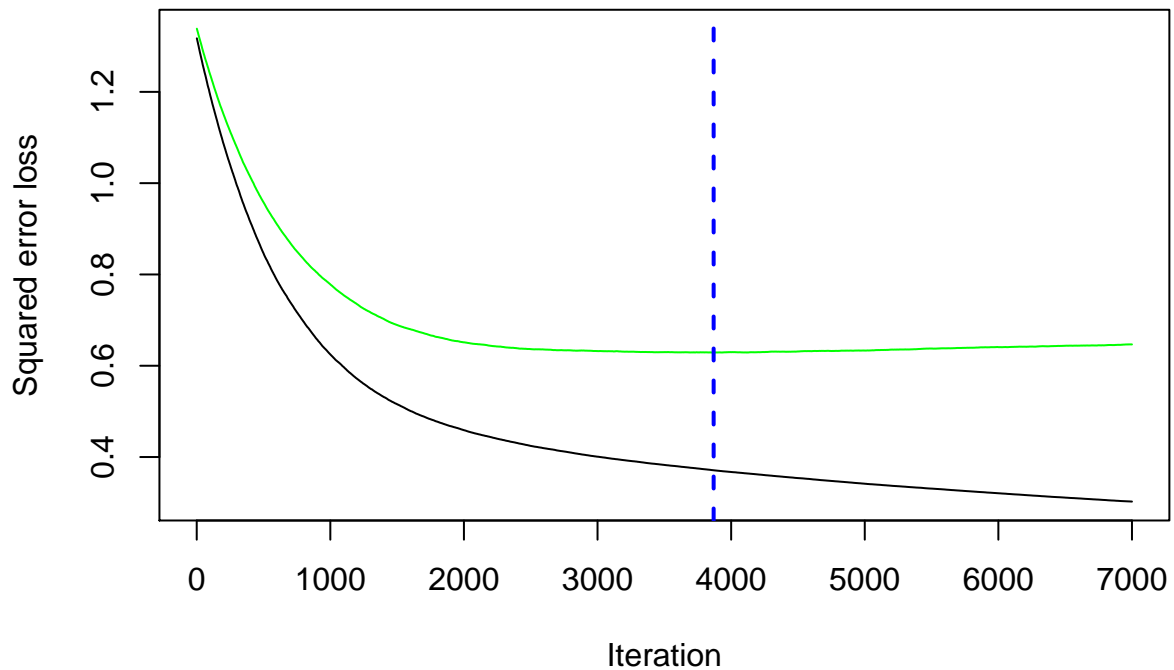
e. Perform boosting and report variable importance

We first fit a gradient boosting model using Gaussian loss function and 7000 iterations. This model only tunes over B (number of trees/iterations) given we have specified fixed interaction depth (d) and shrinkage (λ), although the latter two parameters can also be tuned to select optimal model.

```
set.seed(13)
bst <- gbm(lpsa ~., Prostate,
  distribution = "gaussian",
  n.trees = 7000,
  interaction.depth = 3,
  shrinkage = 0.001,
  cv.folds = 10)
```

We then plot loss function as a result of number of trees added to the ensemble

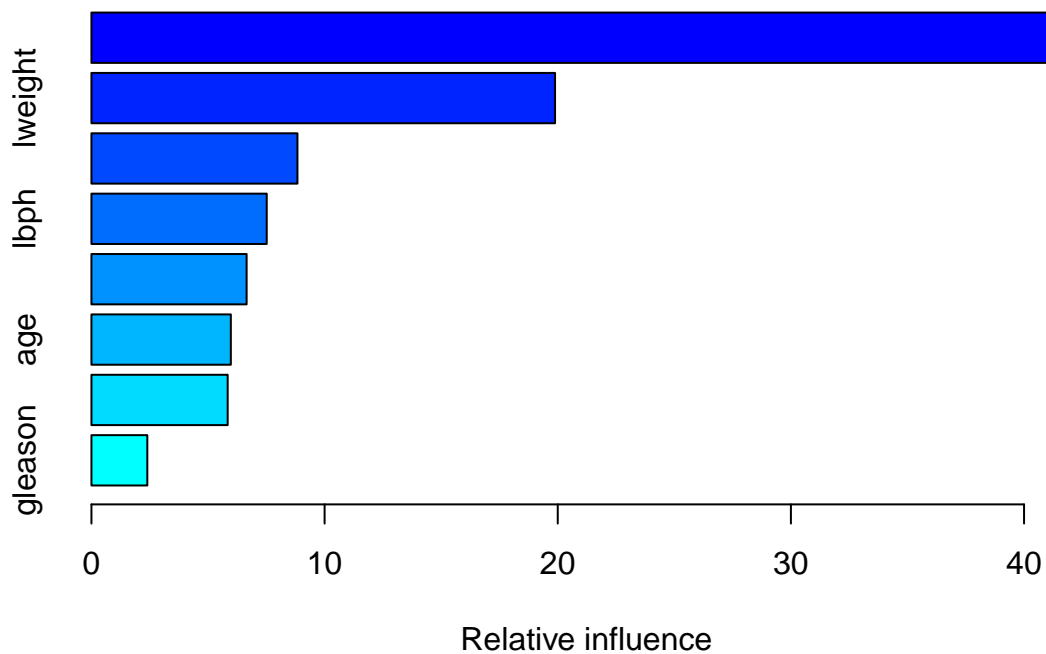
```
nt = gbm.perf(bst, method = "cv")
```



We can see that after 3869 trees, the cross-validated error for the RF model increases. So the optimal number of tree is 3869 for this RF model, given depth = 3, and shrinkage = 0.005.

Obtain variable importance for GBM

```
summary(bst)
```



```
##      var  rel.inf
## lcavol  lcavol 42.888088
## lweight lweight 19.885056
## svi     svi   8.835254
## lbph    lbph  7.518372
## pgg45   pgg45 6.655528
```

```
## age          age  5.978024
## lcp          lcp  5.843455
## gleason      gleason 2.396224
```

From the variable importance matrix and plot we can see log(cancer volume) and log(prostate weight) have the most influence in the fitted models. The plot also suggests age, log(capsular penetration), and gleason carry the least importance in the fitted models. This result is consistent with the result regarding variable importance obtained above.

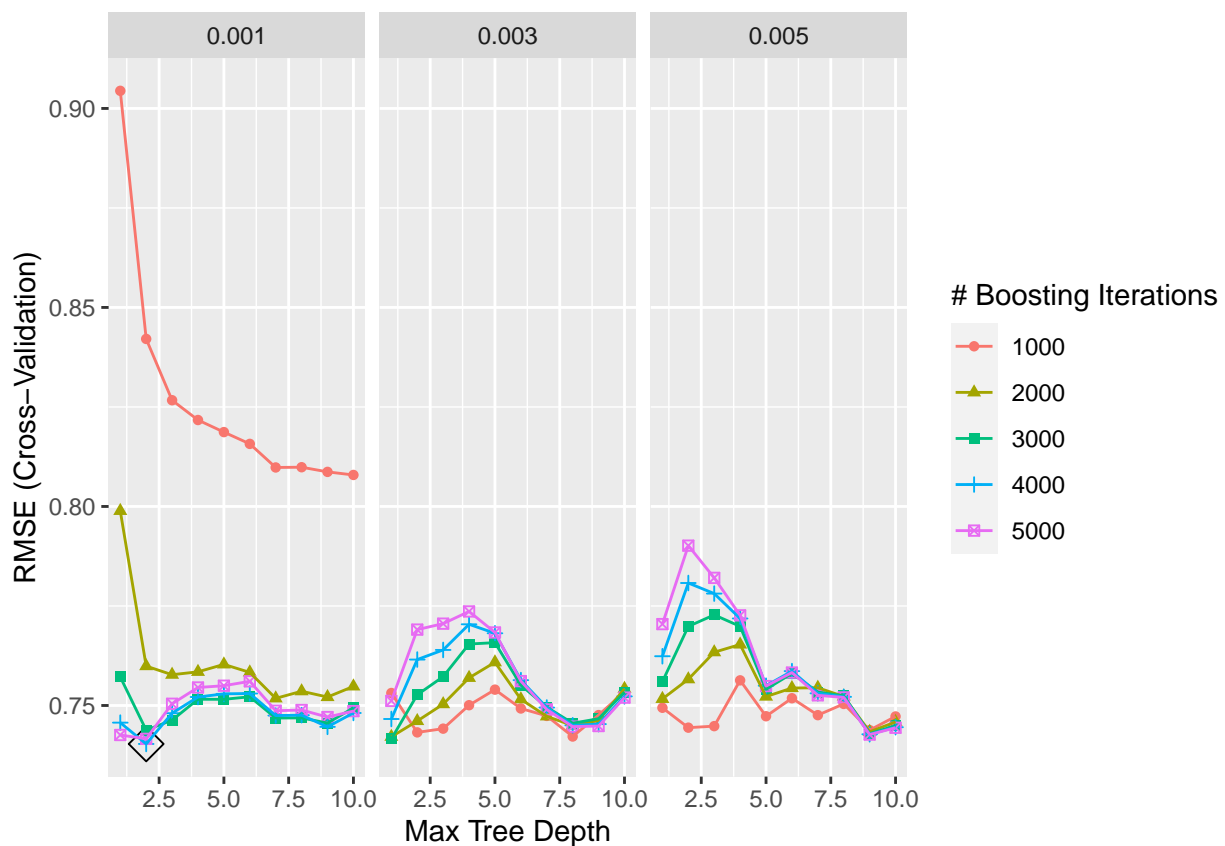
We next try doing a grid search by caret for Boosting, tuning over all 3 parameters

```
gbm.grid = expand.grid(n.trees = seq(1000, 5000, 1000),
                      interaction.depth = 1:10,
                      shrinkage = c(0.001, 0.003, 0.005),
                      n.minobsinnode = 1)
```

```
set.seed(13)
```

```
gbm.fit = train(lpsa~., Prostate,
               method = "gbm",
               tuneGrid = gbm.grid,
               trControl = ctrl,
               verbose = FALSE)
```

```
ggplot(gbm.fit, highlight = TRUE)
```



```
gbm.fit$finalModel$shrinkage
```

```
## [1] 0.001
```



```
gbm.fit$finalModel$interaction.depth
```

```
## [1] 2
```

f. Final choice of model for PSA level prediction

```
resamp <- resamples(list(rf = rf.fit, gbm = gbm.fit))
summary(resamp)
```

```
##
## Call:
## summary.resamples(object = resamp)
##
## Models: rf, gbm
## Number of resamples: 10
##
## MAE
##      Min.    1st Qu.    Median      Mean   3rd Qu.    Max. NA's
## rf  0.2763320 0.5167654 0.6595270 0.6277933 0.7429975 0.9013216    0
## gbm 0.2388993 0.4991992 0.6654952 0.6163851 0.7392922 0.8348874    0
##
## RMSE
##      Min.    1st Qu.    Median      Mean   3rd Qu.    Max. NA's
## rf  0.3562610 0.6050466 0.7655069 0.7590245 0.8833332 1.082792    0
## gbm 0.3138184 0.5862092 0.7741247 0.7403355 0.8947361 1.040386    0
##
## Rsquared
##      Min.    1st Qu.    Median      Mean   3rd Qu.    Max. NA's
## rf  0.4322273 0.5516023 0.6594886 0.6603324 0.7788504 0.8682204    0
## gbm 0.4144193 0.5721892 0.6804740 0.6731744 0.7595333 0.9127964    0
```

From the table, we can see the model obtained from boosting has slightly lower mean RMSE and slightly higher mean R-squared. However, in the end, I decide to go with the model obtained from gradient boosting, with trees = 4000, depth = 2, and shrinkage = 0.001.

Problem 2

Import data (use data OJ in the ISLR package)

```
data("OJ")
oj_data = OJ %>% janitor::clean_names() %>%
  mutate(purchase = as.factor(purchase))
```

a. Fit a classification tree to dataset

```
set.seed(77)
rowTrain = createDataPartition(y = oj_data$purchase,
                                p = 2/3,
                                list = FALSE)

oj_train <- oj_data[rowTrain,]
oj_test  <- oj_data[-rowTrain,]
```

Using “rpart”

```
set.seed(77)
cl_tree = rpart(formula = purchase ~ ., data = oj_train,
```

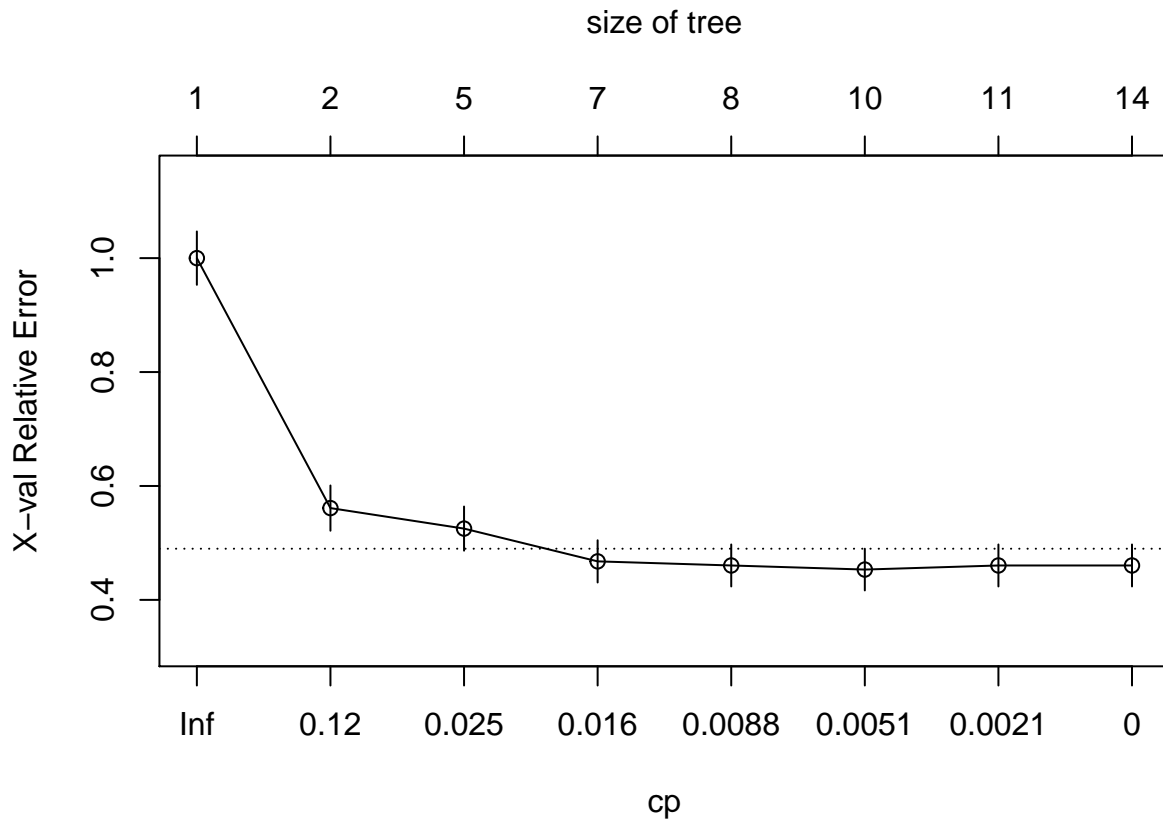
```

control = rpart.control(cp = 0))

cpTable = printcp(cl_tree)

##
## Classification tree:
## rpart(formula = purchase ~ ., data = oj_train, control = rpart.control(cp = 0))
##
## Variables actually used in tree construction:
## [1] list_price_diff loyal_ch      price_diff      special_ch
## [5] store          weekof_purchase
##
## Root node error: 278/714 = 0.38936
##
## n= 714
##
##      CP nsplit rel error  xerror   xstd
## 1 0.5000000    0  1.00000 1.00000 0.046867
## 2 0.0275779    1  0.50000 0.56115 0.039718
## 3 0.0233813    4  0.41727 0.52518 0.038767
## 4 0.0107914    6  0.37050 0.46763 0.037092
## 5 0.0071942    7  0.35971 0.46043 0.036869
## 6 0.0035971    9  0.34532 0.45324 0.036642
## 7 0.0011990   10  0.34173 0.46043 0.036869
## 8 0.0000000   13  0.33813 0.46043 0.036869
plotcp(cl_tree)

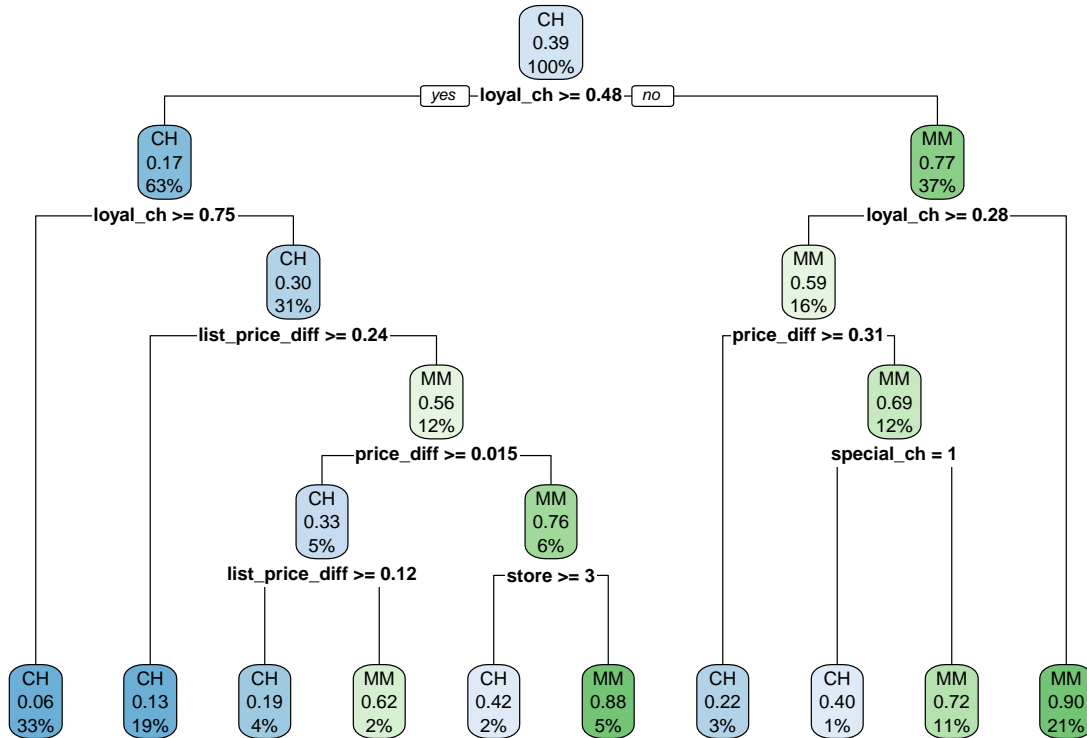
```



Use cross-validation to determine the tree size and create a plot of the final tree.

```
# min CV error
minErr = which.min(cpTable[,4])

# pruning
cl_tree1 = prune(cl_tree, cp = cpTable[minErr,1])
rpart.plot(cl_tree1)
```



Predict the response on the test data and find the test error rate.

```
tree_pred <- predict(cl_tree1, oj_test, type = "class")
head(tree_pred, 10)
```

```
## 1 4 6 7 8 9 10 15 17 19
## CH MM CH CH CH CH CH MM CH CH
## Levels: CH MM
```

The predicted class labels for the first 10 observations in the test data can be observed as above.

```
table(tree_pred, oj_test$purchase)
```

```
##
## tree_pred CH MM
##          CH 187 37
##          MM 30 102
```

From the confusion matrix, we can calculate the test error rate as: $(30+37)/(187+37+30+102) = 0.19$.

b. Perform random forests on the training set and report variable importance. Find test error rate

```
set.seed(77)
rf = ranger(purchase~., data = oj_train,
```

```

mtry = 6, probability = TRUE)

rf.pred = predict(rf, oj_test, type = "response")$predictions[,1]
head(rf.pred, 20)

```

```

## [1] 0.5882220 0.0711364 0.9301036 0.9474025 0.8083132 0.9239975 0.9226934
## [8] 0.7440718 0.5922476 0.6993254 0.9868731 0.3223393 0.3679458 0.9799079
## [15] 0.9186659 0.4733796 0.8794714 0.9085238 0.9237571 0.7149212

```

We can see the first 20 predictions in the test dataset. Observations with probability < 0.5 will be classified into “Minute Maid”, while with probability > 0.5 will be classified into “Citrus Hill.”

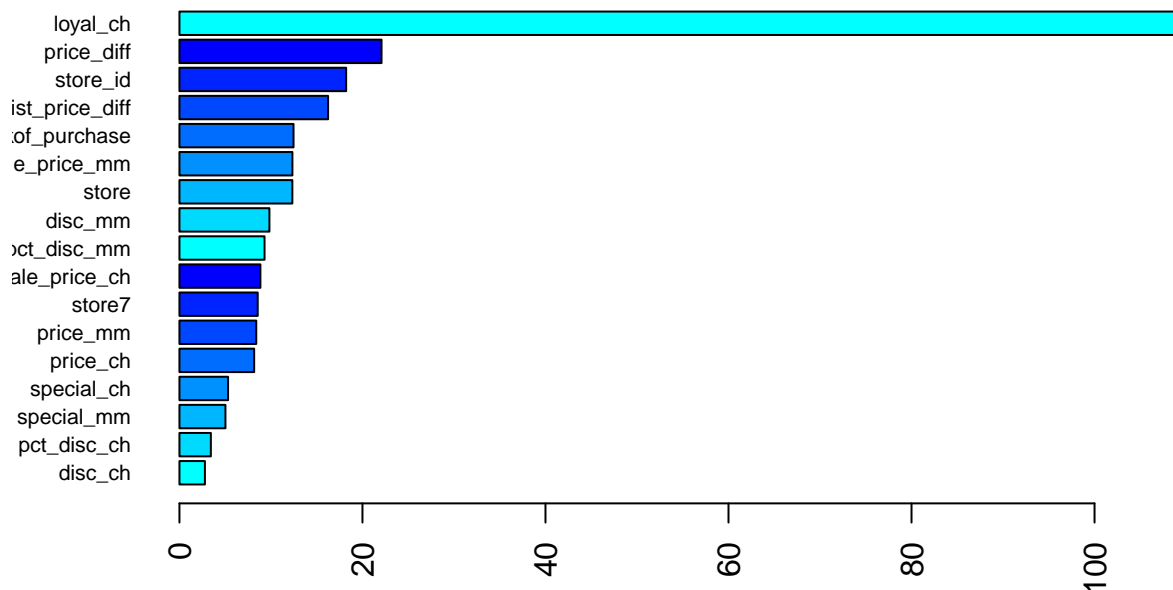
Obtain variable importance

```

set.seed(77)
rf2.final.per <- ranger(purchase~., oj_train,
                        mtry = 6, min.node.size = 5,
                        splitrule = "gini",
                        importance = "permutation",
                        scale.permutation.importance = TRUE)

barplot(sort(ranger::importance(rf2.final.per), decreasing = FALSE),
        las = 2, horiz = TRUE, cex.names = 0.7, col = colorRampPalette(colors = c("cyan", "blue"))(8))

```



We can see loyal_ch (customer brand loyalty for Citrus Hill), price_diff (sale price of MM minus sale price of CH), and store ID are the most influential features, whereas percent discount and discount offered for Citrus Hill are the least important in the random forest fitted model.

```

rf.pred = ifelse(rf.pred > 0.5, "CH", "MM")
table(rf.pred, oj_test$purchase)

```

```

##
## rf.pred  CH  MM
##      CH 187  43
##      MM  30  96

```

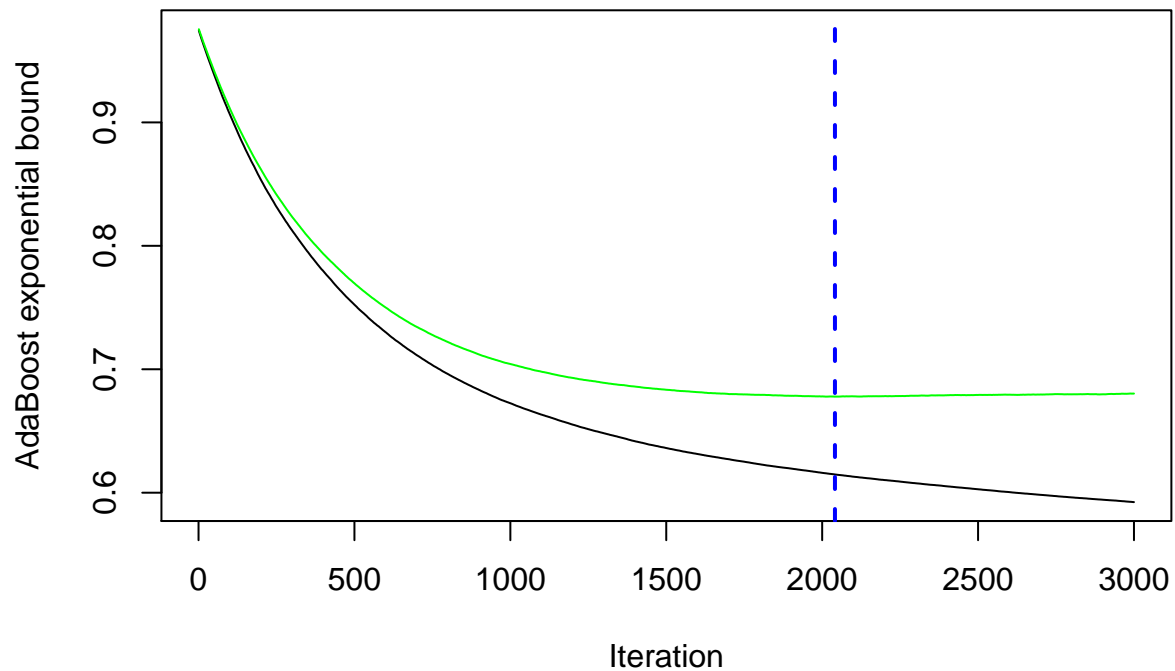
The test error rate can be calculated as $(30+43)/(187+43+30+96) = 0.21$

c. Perform boosting and report variable importance. Obtain test error rate

```
oj_train$purchase <-as.numeric(oj_train$purchase=="MM")
oj_test$purchase <-as.numeric(oj_test$purchase=="MM")
```

```
set.seed(77)
bst1 <-gbm(purchase~., oj_train,
  distribution = "adaboost",
  n.trees = 3000,interaction.depth = 2,
  shrinkage = 0.002,
  cv.folds = 10)

nt1 <-gbm.perf(bst1, method = "cv")
```

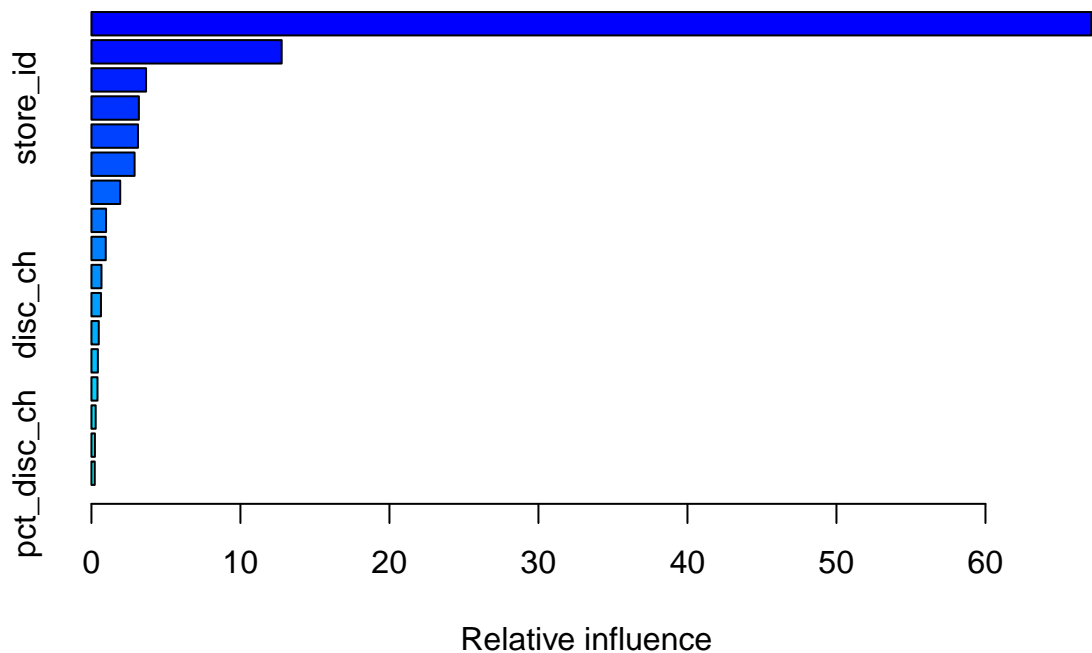


```
nt1
```

```
## [1] 2041
```

The number of optimal trees/iterations picked by gbm function is 2041, which minimizes the error loss given interaction depth = 3 and shrinkage = 0.005.

```
summary(bst1)
```



```
##           var      rel.inf
## loyal_ch      loyal_ch 67.1094987
## price_diff    price_diff 12.7678944
## weekof_purchase weekof_purchase 3.6661785
## store_id      store_id 3.1836302
## list_price_diff list_price_diff 3.1234904
## sale_price_mm  sale_price_mm 2.8956804
## store         store 1.9303689
## sale_price_ch  sale_price_ch 0.9811661
## price_mm      price_mm 0.9602485
## pct_disc_mm   pct_disc_mm 0.6748082
## disc_ch       disc_ch 0.6385541
## disc_mm       disc_mm 0.4914635
## store7        store7 0.4315498
## special_mm    special_mm 0.4064008
## price_ch      price_ch 0.2843065
## special_ch    special_ch 0.2329142
## pct_disc_ch   pct_disc_ch 0.2218467
```

We can observe similar results as in RF model regarding variable importance, e.g., “loyal_ch” (customer brand loyalty for Citrus Hill), price_diff (sale price of MM minus sale price of CH) and store id are still the most influential features, whereas percent discount and discount offered for Citrus Hill are still among the least important variables.

Obtain test error rate

```
gbm.pred = predict(bst1, newdata = oj_test, type = "response")
```

```
## Using 2041 trees...
```

```
head(gbm.pred, 20)
```

```
## [1] 0.43654903 0.76034684 0.03923052 0.06378494 0.06837687 0.06865853
## [7] 0.06733617 0.12179071 0.55082177 0.48784798 0.03567037 0.53423840
## [13] 0.66647244 0.07401541 0.04042834 0.67841139 0.11231166 0.11092811
```

```
## [19] 0.10771026 0.05215241
```

Similarly, we can see the first 20 predictions in the test dataset. Observations with probability < 0.5 will be classified into “Minute Maid”, while with probability > 0.5 will be classified into “Citrus Hill.”

```
gbm.pred = ifelse(gbm.pred > 0.5, "CH", "MM")
table(gbm.pred, oj_test$purchase)
```

```
##
## gbm.pred    0    1
##          CH  27 102
##          MM 190  37
```

With $CH = 1$, $MM = 0$, the test error rate for this model is $(27+37)/(27+190+102+37) = 0.18$