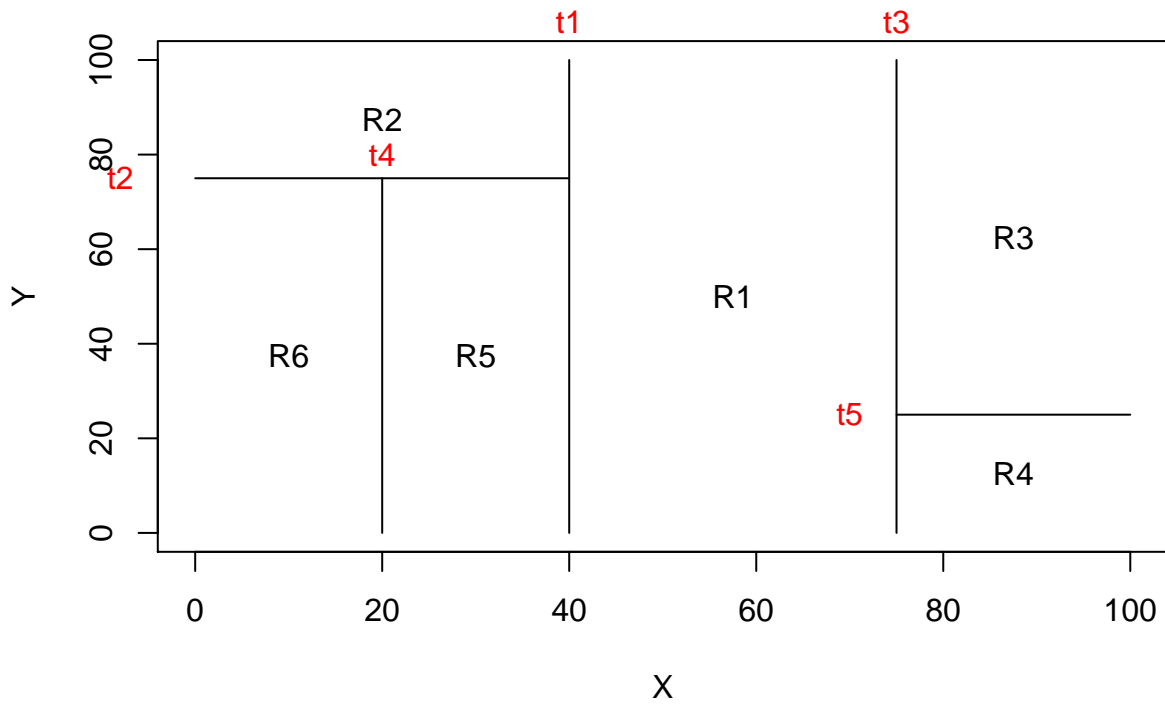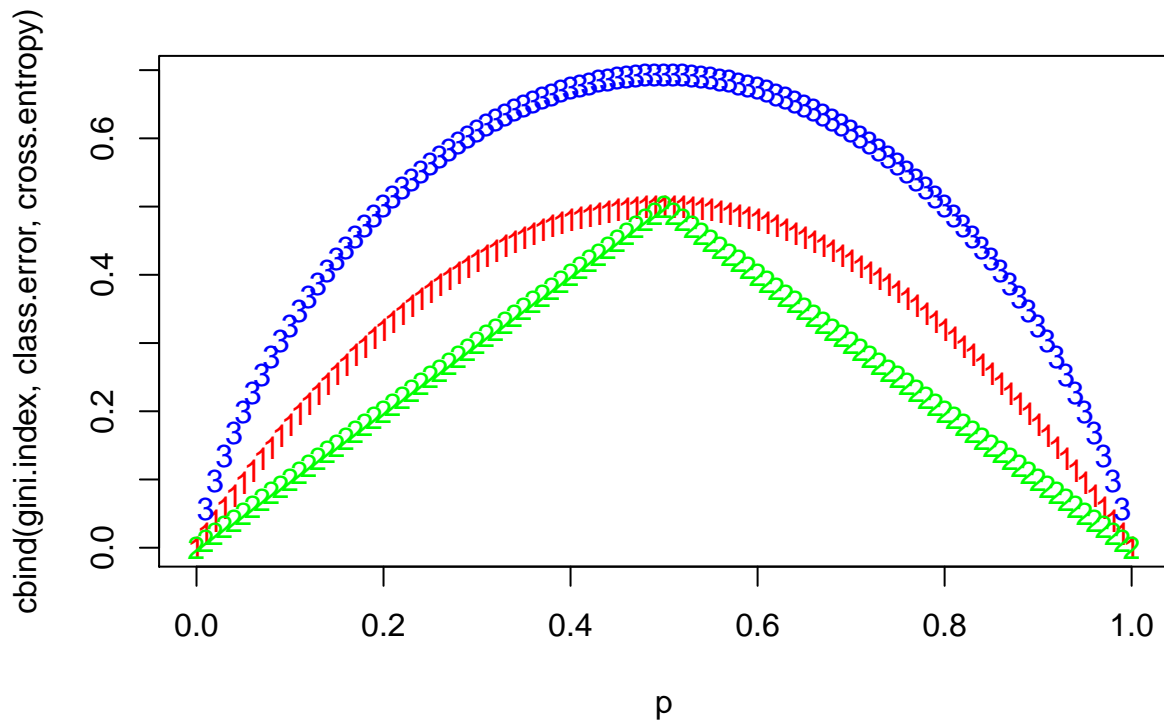# chp8 tree

## Wendy Liang

## 2/28/2021

**8.1**

```r
par(xpd = NA)
plot(NA, NA, type = "n", xlim = c(0,100), ylim = c(0,100), xlab = "X", ylab = "Y")
# t1: x = 40; (40, 0) (40, 100)
lines(x = c(40,40), y = c(0,100))
text(x = 40, y = 108, labels = c("t1"), col = "red")
# t2: y = 75; (0, 75) (40, 75)
lines(x = c(0,40), y = c(75,75))
text(x = -8, y = 75, labels = c("t2"), col = "red")
# t3: x = 75; (75,0) (75, 100)
lines(x = c(75,75), y = c(0,100))
text(x = 75, y = 108, labels = c("t3"), col = "red")
# t4: x = 20; (20,0) (20, 75)
lines(x = c(20,20), y = c(0,75))
text(x = 20, y = 80, labels = c("t4"), col = "red")
# t5: y=25; (75,25) (100,25)
lines(x = c(75,100), y = c(25,25))
text(x = 70, y = 25, labels = c("t5"), col = "red")

text(x = (40+75)/2, y = 50, labels = c("R1"))
text(x = 20, y = (100+75)/2, labels = c("R2"))
text(x = (75+100)/2, y = (100+25)/2, labels = c("R3"))
text(x = (75+100)/2, y = 25/2, labels = c("R4"))
text(x = 30, y = 75/2, labels = c("R5"))
text(x = 10, y = 75/2, labels = c("R6"))
```

**8.3**

```r
p <- seq(0, 1, 0.01)
gini.index <- 2 * p * (1 - p)
class.error <- 1 - pmax(p, 1 - p)
cross.entropy <- - (p * log(p) + (1 - p) * log(1 - p))
matplot(p, cbind(gini.index, class.error, cross.entropy), col = c("red", "green", "blue"))
```
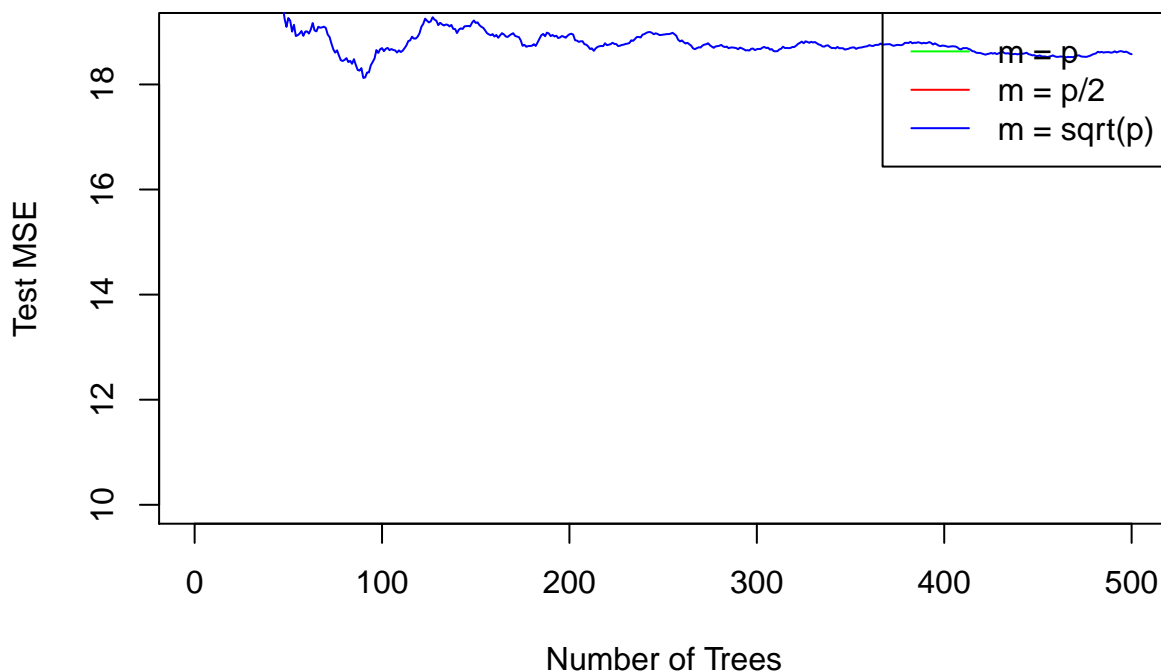
**8.5**

With the majority vote approach, we classify X as Red as it is the most commonly occurring class among the 10 predictions (6 for Red vs 4 for Green). With the average probability approach, we classify X as Green as the average of the 10 probabilities is 0.45.

**8.7**

```
set.seed(1)
train <- sample(1:nrow(Boston), nrow(Boston) / 2)
Boston.train <- Boston[train, -14]
Boston.test <- Boston[-train, -14]
Y.train <- Boston[train, 14]
Y.test <- Boston[-train, 14]
rf.boston1 <- randomForest(Boston.train, y = Y.train, xtest = Boston.test, ytest = Y.test, mtry = ncol(
rf.boston2 <- randomForest(Boston.train, y = Y.train, xtest = Boston.test, ytest = Y.test, mtry = (ncol
rf.boston3 <- randomForest(Boston.train, y = Y.train, xtest = Boston.test, ytest = Y.test, mtry = sqrt(
plot(1:500, rf.boston1$test$mse, col = "green", type = "l", xlab = "Number of Trees", ylab = "Test MSE"
lines(1:500, rf.boston2$test$mse, col = "red", type = "l")
lines(1:500, rf.boston3$test$mse, col = "blue", type = "l")
legend("topright", c("m = p", "m = p/2", "m = sqrt(p)"), col = c("green", "red", "blue"), cex = 1, lty =
```



We may see that the Test MSE is very high for a single tree, it decreases as the number of trees increases. Also the Test MSE for all predictors is higher than for half the predictors or the square root of the number of predictors.
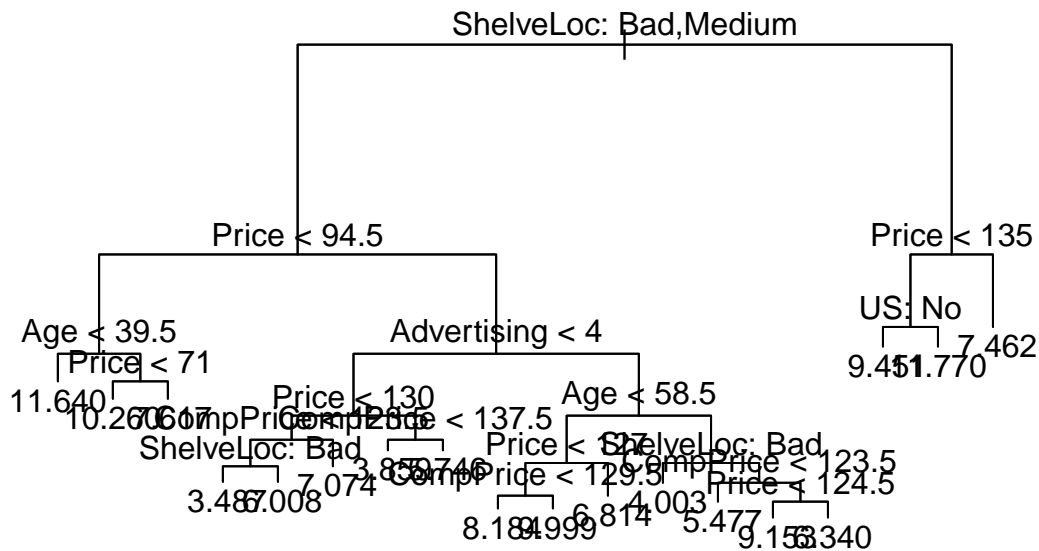
**8.8**

```
set.seed(1)
train <- sample(1:nrow(Carseats), nrow(Carseats) / 2)
Carseats.train <- Carseats[train, ]
Carseats.test <- Carseats[-train, ]
```

**a**

```
tree.carseats <- tree(Sales ~ ., data = Carseats.train)
summary(tree.carseats)
```

**b**

```
##
## Regression tree:
## tree(formula = Sales ~ ., data = Carseats.train)
## Variables actually used in tree construction:
## [1] "ShelveLoc"   "Price"        "Age"          "Advertising" "CompPrice"
## [6] "US"
## Number of terminal nodes:  18
## Residual mean deviance:  2.167 = 394.3 / 182
## Distribution of residuals:
##     Min.  1st Qu.   Median     Mean  3rd Qu.     Max.
## -3.88200 -0.88200 -0.08712  0.00000  0.89590  4.09900
```

```
plot(tree.carseats)
text(tree.carseats, pretty = 0)
```
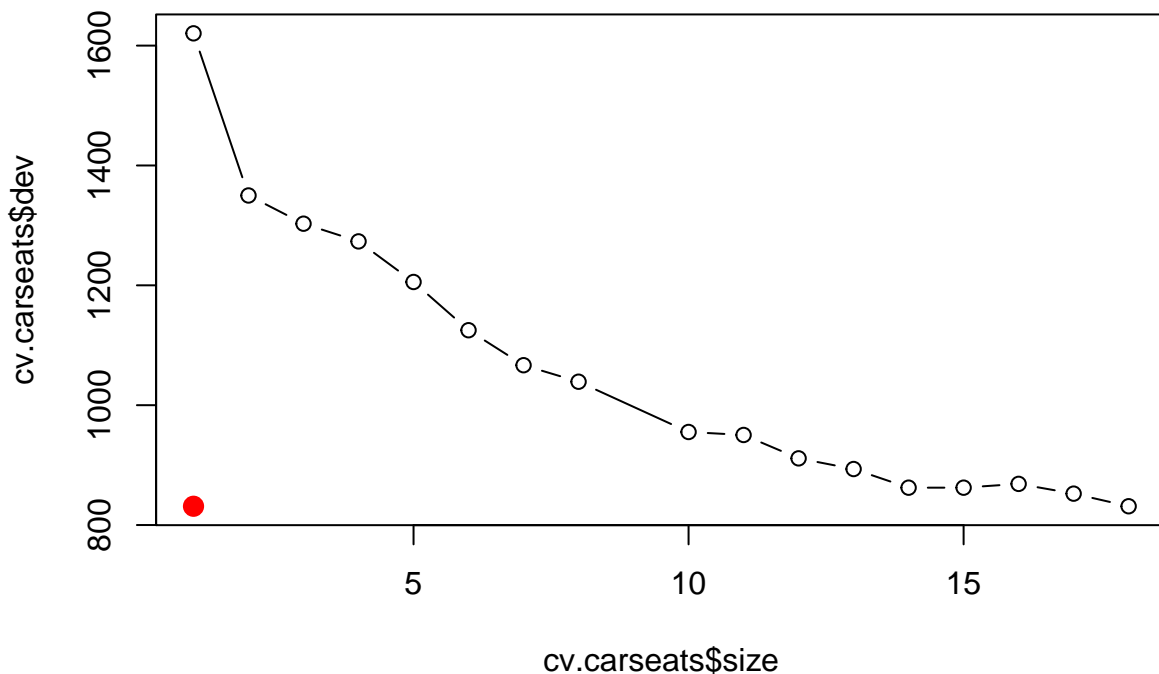


```
yhat <- predict(tree.carseats, newdata = Carseats.test)
mean((yhat - Carseats.test$Sales)^2)
```
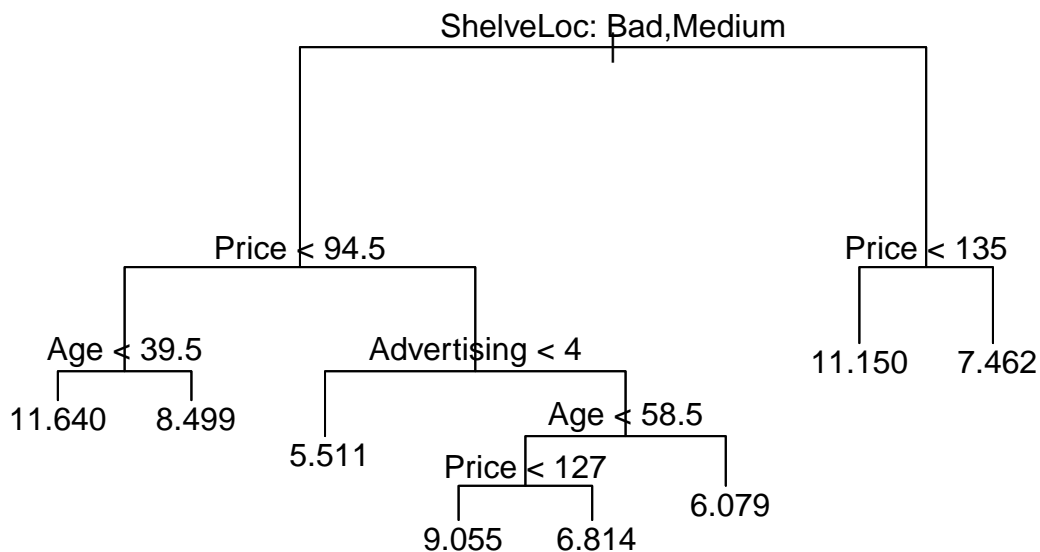
```
## [1] 4.922039
```

```
cv.carseats <- cv.tree(tree.carseats)
plot(cv.carseats$size, cv.carseats$dev, type = "b")
tree.min <- which.min(cv.carseats$dev)
points(tree.min, cv.carseats$dev[tree.min], col = "red", cex = 2, pch = 20)
```

**c**

```
#In this case, the tree of size 8 is selected by cross-validation.
#We now prune the tree to obtain the 8-node tree.

prune.carseats <- prune.tree(tree.carseats, best = 8)
plot(prune.carseats)
text(prune.carseats, pretty = 0)
```



```
yhat <- predict(prune.carseats, newdata = Carseats.test)
mean((yhat - Carseats.test$Sales)^2)
```

```
## [1] 5.113254
```

```
bag.carseats <- randomForest(Sales ~ ., data = Carseats.train, mtry = 10, ntree = 500, importance = TRUE
yhat.bag <- predict(bag.carseats, newdata = Carseats.test)
```

```r
mean((yhat.bag - Carseats.test$Sales)^2)
```

**d**

```
## [1] 2.657296
```

```r
importance(bag.carseats)
```

```
##                 %IncMSE IncNodePurity
## CompPrice     23.07909904    171.185734
## Income         2.82081527     94.079825
## Advertising   11.43295625     99.098941
## Population    -3.92119532     59.818905
## Price         54.24314632    505.887016
## ShelveLoc     46.26912996    361.962753
## Age           14.24992212    159.740422
## Education     -0.07662320     46.738585
## Urban          0.08530119      8.453749
## US             4.34349223     15.157608
```

```r
rf.carseats <- randomForest(Sales ~ ., data = Carseats.train, mtry = 3, ntree = 500, importance = TRUE)
yhat.rf <- predict(rf.carseats, newdata = Carseats.test)
mean((yhat.rf - Carseats.test$Sales)^2)
```

**e**

```
## [1] 3.049406
```

```r
importance(rf.carseats)
```

```
##                %IncMSE IncNodePurity
## CompPrice    12.9489323     158.48521
## Income        2.2754686     129.59400
## Advertising   8.9977589     111.94374
## Population   -2.2513981     102.84599
## Price        33.4226950     391.60804
## ShelveLoc    34.0233545     290.56502
## Age          12.2185108     171.83302
## Education     0.2592124      71.65413
## Urban         1.1382113      14.76798
## US            4.1925335      33.75554
```
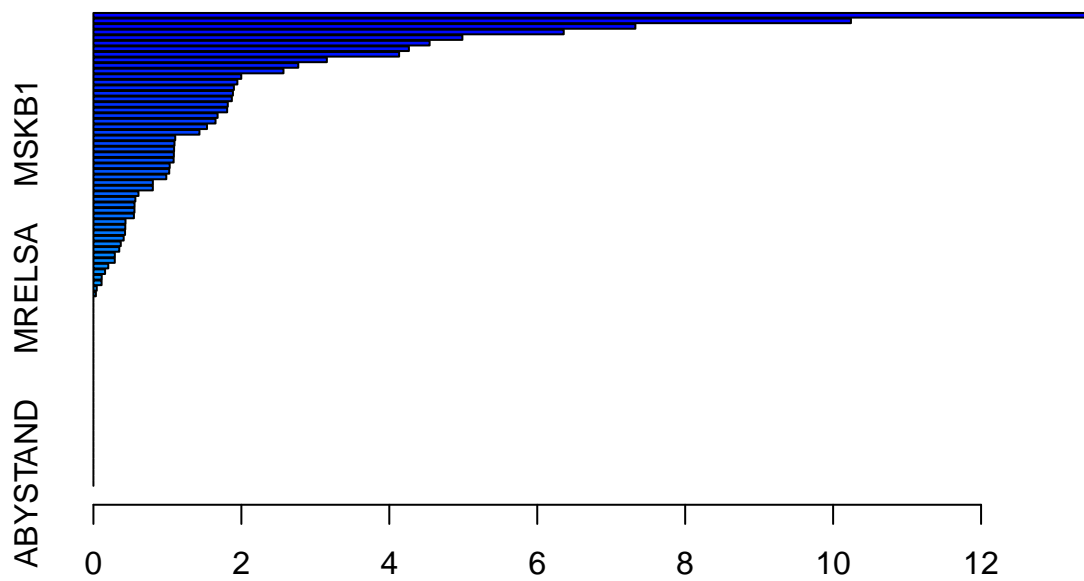
**8.11**

```r
set.seed(1)
train <- 1:1000
Caravan$Purchase <- ifelse(Caravan$Purchase == "Yes", 1, 0)
Caravan.train <- Caravan[train, ]
Caravan.test <- Caravan[-train, ]
```

**a**

```r
set.seed(1)
boost.caravan <- gbm(Purchase ~ ., data = Caravan.train, distribution = "gaussian", n.trees = 1000, shr
```

```
summary(boost.caravan)
```

**b**



Relative influence

```
##                    var      rel.inf
## PPERSAUT PPERSAUT 13.51824557
## MKOOPKLA MKOOPKLA 10.24062778
## MOPLHOOG MOPLHOOG  7.32689780
## MBERMIDD MBERMIDD  6.35820558
## PBRAND     PBRAND   4.98826360
## ABRAND     ABRAND   4.54504653
## MGODGE     MGODGE   4.26496875
## MINK3045 MINK3045  4.13253907
## PWAPART   PWAPART   3.15612877
## MAUT1       MAUT1   2.76929763
## MOSTYPE   MOSTYPE   2.56937935
## MAUT2       MAUT2   1.99879666
## MSKA         MSKA   1.94618539
## MBERARBG MBERARBG  1.89917331
## PBYSTAND PBYSTAND  1.88591514
## MINKGEM   MINKGEM   1.87131472
## MGODOV     MGODOV   1.81673309
## MGODPR     MGODPR   1.80814745
## MFWEKIND MFWEKIND  1.67884570
## MSKC         MSKC   1.65075962
## MBERHOOG MBERHOOG  1.53559951
## MSKB1       MSKB1   1.43339514
## MOPLMIDD MOPLMIDD  1.10617074
## MHHUUR     MHHUUR   1.09608784
## MRELGE     MRELGE   1.09039794
## MINK7512 MINK7512  1.08772012
## MZFONDS   MZFONDS   1.08427551
## MGODRK     MGODRK   1.03126657
```

```
## MINK4575 MINK4575  1.02492795
## MZPART    MZPART    0.98536712
## MRELOV    MRELOV    0.80356854
## MFGEKIND MFGEKIND   0.80335689
## MBERARBO MBERARBO   0.60909852
## APERSAUT APERSAUT   0.56707821
## MGEMOMV   MGEMOMV   0.55589456
## MOSHOOFD MOSHOOFD   0.55498375
## MAUTO      MAUTO    0.54748481
## PMOTSCO   PMOTSCO   0.43362597
## MSKB2      MSKB2    0.43075446
## MSKD        MSKD    0.42751490
## MINK123M MINK123M   0.40920707
## MINKM30   MINKM30   0.36996576
## MHKOOP    MHKOOP    0.34941518
## MBERBOER MBERBOER   0.28967068
## MFALLEEN MFALLEEN   0.28877552
## MGEMLEEF MGEMLEEF   0.20084195
## MOPLLAAG MOPLLAAG   0.15750616
## MBERZELF MBERZELF   0.11203381
## PLEVEN     PLEVEN   0.11030994
## MRELSA    MRELSA    0.04500507
## MAANTHUI MAANTHUI   0.03322830
## PWABEDR   PWABEDR   0.00000000
## PWALAND   PWALAND   0.00000000
## PBESAUT   PBESAUT   0.00000000
## PVRAAUT   PVRAAUT   0.00000000
## PAANHANG PAANHANG   0.00000000
## PTRACTOR PTRACTOR   0.00000000
## PWERKT     PWERKT   0.00000000
## PBROM       PBROM   0.00000000
## PPERSONG PPERSONG   0.00000000
## PGEZONG   PGEZONG   0.00000000
## PWAOREG   PWAOREG   0.00000000
## PZEILPL   PZEILPL   0.00000000
## PPLEZIER PPLEZIER   0.00000000
## PFIETS     PFIETS   0.00000000
## PINBOED   PINBOED   0.00000000
## AWAPART   AWAPART   0.00000000
## AWABEDR   AWABEDR   0.00000000
## AWALAND   AWALAND   0.00000000
## ABESAUT   ABESAUT   0.00000000
## AMOTSCO   AMOTSCO   0.00000000
## AVRAAUT   AVRAAUT   0.00000000
## AAANHANG AAANHANG   0.00000000
## ATRACTOR ATRACTOR   0.00000000
## AWERKT     AWERKT   0.00000000
## ABROM       ABROM   0.00000000
## ALEVEN     ALEVEN   0.00000000
## APERSONG APERSONG   0.00000000
## AGEZONG   AGEZONG   0.00000000
## AWAOREG   AWAOREG   0.00000000
## AZEILPL   AZEILPL   0.00000000
## APLEZIER APLEZIER   0.00000000
```

```
## AFIETS      AFIETS   0.00000000
## AINBOED    AINBOED   0.00000000
## ABYSTAND ABYSTAND   0.00000000
```

```
probs.test <- predict(boost.caravan, Caravan.test, n.trees = 1000, type = "response")
pred.test <- ifelse(probs.test > 0.2, 1, 0)
table(Caravan.test$Purchase, pred.test)
```

c

```
##    pred.test
##        0    1
##   0 4493   40
##   1  278   11
```

```
logit.caravan <- glm(Purchase ~ ., data = Caravan.train, family = "binomial")

probs.test2 <- predict(logit.caravan, Caravan.test, type = "response")

pred.test2 <- ifelse(probs.test > 0.2, 1, 0)
table(Caravan.test$Purchase, pred.test2)
```

```
##    pred.test2
##        0    1
##   0 4493   40
##   1  278   11
```