

R 실습

나무모형

임요한

서울대학교

Sep, 2018

패키지 로딩

```
library(ISLR)  
attach(Wage)
```

1. 분류나무

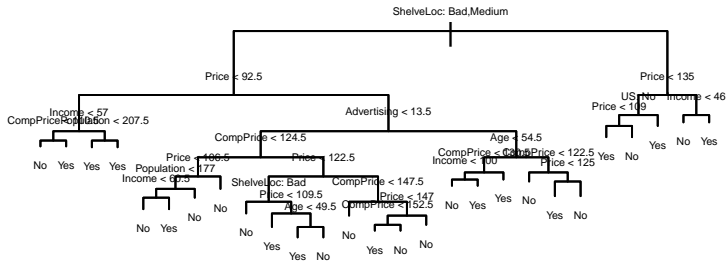
Fitting Classification Trees

```
library(ISLR)
library(tree)
attach(Carseats)
High=ifelse(Sales<=8,"No","Yes")
Carseats=data.frame(Carseats,High)
```

```
tree.carseats=tree(High ~ . - Sales, Carseats)
summary(tree.carseats)
```

```
##
## Classification tree:
## tree(formula = High ~ . - Sales, data = Carseats)
## Variables actually used in tree construction:
## [1] "ShelveLoc"    "Price"        "Income"       "CompPrice"    "Popu
## [6] "Advertising" "Age"          "US"
## Number of terminal nodes: 27
## Residual mean deviance: 0.4575 = 170.7 / 373
## Misclassification error rate: 0.09 = 36 / 400
```

```
plot(tree.carseats)
text(tree.carseats, pretty=0, cex=0.3)
```



```
## node), split, n, deviance, yval, (yprob)
```

```
##      * denotes terminal node
```

```
##
```

```
## 1) root 400 541.500 No ( 0.59000 0.41000 )
```

```
## 2) ShelfLoc: Bad,Medium 315 390.600 No ( 0.68889 0.31111 )
```

```
## 4) Price < 92.5 46 56.530 Yes ( 0.30435 0.69565 )
```

```
## 8) Income < 57 10 12.220 No ( 0.70000 0.30000 )
```

```
## 16) CompPrice < 110.5 5 0.000 No ( 1.00000 0.00000 ) *
```

```
## 17) CompPrice > 110.5 5 6.730 Yes ( 0.40000 0.60000 )
```

```
## 9) Income > 57 36 35.470 Yes ( 0.19444 0.80556 )
```

```
## 18) Population < 207.5 16 21.170 Yes ( 0.37500 0.62500 )
```

```
## 19) Population > 207.5 20 7.941 Yes ( 0.05000 0.95000 )
```

```
## 5) Price > 92.5 269 299.800 No ( 0.75465 0.24535 )
```

```
## 10) Advertising < 13.5 224 213.200 No ( 0.81696 0.18304 )
```

```
## 20) CompPrice < 124.5 96 44.890 No ( 0.93750 0.06250 )
```

```
## 40) Price < 106.5 38 33.150 No ( 0.84211 0.15789 )
```

```
## 80) Population < 177 12 16.300 No ( 0.58333 0.41667 )
```

```
## 160) Income < 60.5 6 0.000 No ( 1.00000 0.00000 )
```

```
## 161) Income > 60.5 6 5.407 Yes ( 0.16667 0.83333 )
```

```

set.seed(2)
train=sample(1:nrow(Carseats), 200)
Carseats.test=Carseats[-train,]
High.test=High[-train]
tree.carseats=tree(High~.-Sales,Carseats,subset=train)
tree.pred=predict(tree.carseats,Carseats.test,type="class")
table(tree.pred,High.test)

```

```

##           High.test
## tree.pred No  Yes
##           No  86  27
##           Yes 30  57

```

```

(86+57)/200

```

```

## [1] 0.715    30%  error가

```

```
set.seed(3)
cv.carseats=cv.tree(tree.carseats,FUN=prune.misclass)
names(cv.carseats)
```

```
## [1] "size"    "dev"     "k"       "method"
```

```
cv.carseats
```

```
## $size
```

```
## [1] 19 17 14 13 9 7 3 2 1
```

```
##
```

```
## $dev
```

```
## [1] 55 55 53 52 50 56 69 65 80
```

```
##
```

```
## $k
```

```
## [1] -Inf 0.0000000 0.6666667 1.0000000 1.7500000 2.0000000
```

```
## [7] 4.2500000 5.0000000 23.0000000
```

```
##
```

```
## $method
```

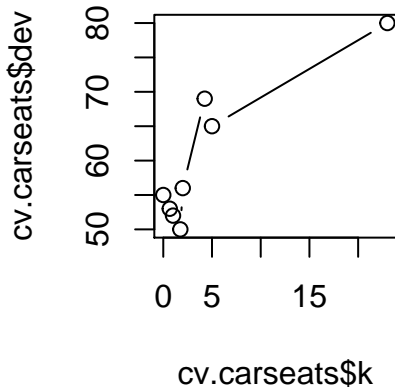
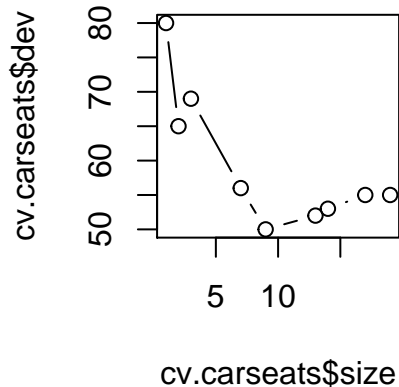
```
## [1] "misclass"
```

```
##
```

```
## attr( "class")
```



```
par(mfrow=c(1,2))  
plot(cv.carseats$size,cv.carseats$dev,type="b")  
plot(cv.carseats$k,cv.carseats$dev,type="b")
```



```
par(mfrow=c(1,1))
prune.carseats=prune.misclass(tree.carseats,best=9)
plot(prune.carseats)
text(prune.carseats,pretty=0,cex=0.5)
```



```
tree.pred=predict(prune.carseats,Carseats.test,type="class")
table(tree.pred,High.test)
```

```
##           High.test
## tree.pred No  Yes
##           No  94  24
##           Yes 22  60
```

```
(94+60)/200
```

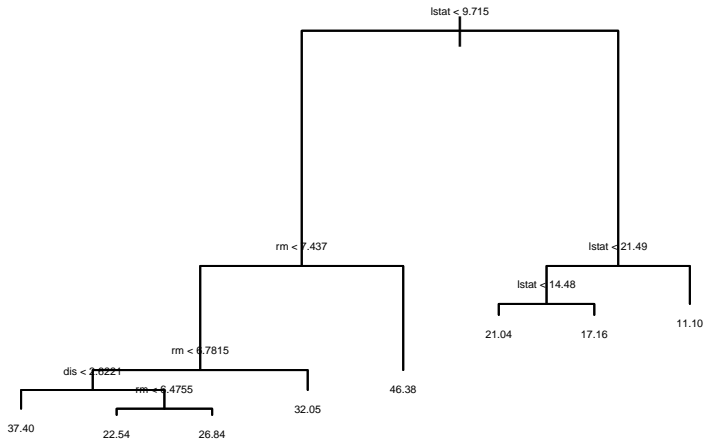
```
## [1] 0.77
```

2. 회귀나무

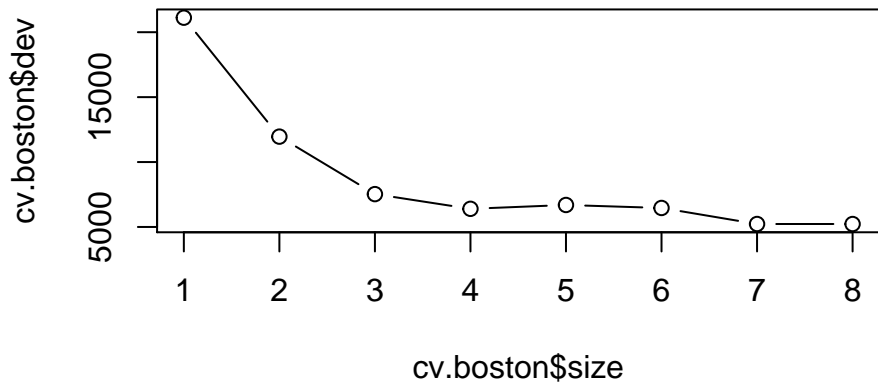
```
library(MASS)
set.seed(1)
train = sample(1:nrow(Boston), nrow(Boston)/2)
tree.boston=tree(medv~.,Boston,subset=train)
summary(tree.boston)
```

```
##
## Regression tree:
## tree(formula = medv ~ ., data = Boston, subset = train)
## Variables actually used in tree construction:
## [1] "lstat" "rm"      "dis"
## Number of terminal nodes: 8
## Residual mean deviance: 12.65 = 3099 / 245
## Distribution of residuals:
##      Min.    1st Qu.    Median      Mean    3rd Qu.      Max.
## -14.10000  -2.04200  -0.05357   0.00000   1.96000  12.60000
```

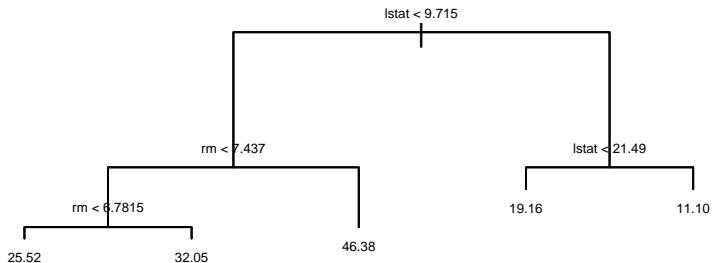
```
plot(tree.boston)
text(tree.boston,pretty=0,cex=0.3)
```



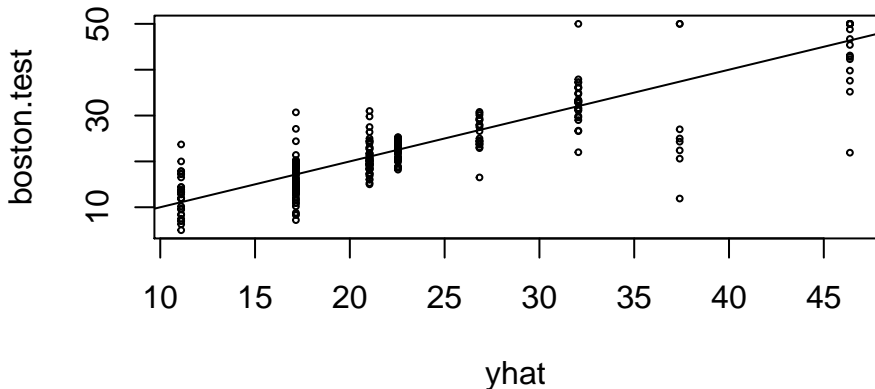
```
cv.boston=cv.tree(tree.boston)
plot(cv.boston$size,cv.boston$dev,type='b')
```



```
prune.boston=prune.tree(tree.boston,best=5)
plot(prune.boston)
text(prune.boston,pretty=0,cex=0.4)
```



```
yhat=predict(tree.boston,newdata=Boston[-train,])  
boston.test=Boston[-train,"medv"]  
plot(yhat,boston.test,cex=0.4)  
abline(0,1)
```



```
mean((yhat-boston.test)^2)
```

```
## [1] 25.04559
```


3. 배깅과 나무숲

```
library(randomForest)
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
set.seed(1)
```

```
bag.boston=randomForest(medv~.,data=Boston,subset=train,mtry=13,importance=TRUE)  
bag.boston
```

```
##
```

```
## Call:
```

```
## randomForest(formula = medv ~ ., data = Boston, mtry = 13, impor
```

```
##           Type of random forest: regression
```

```
##           Number of trees: 500
```

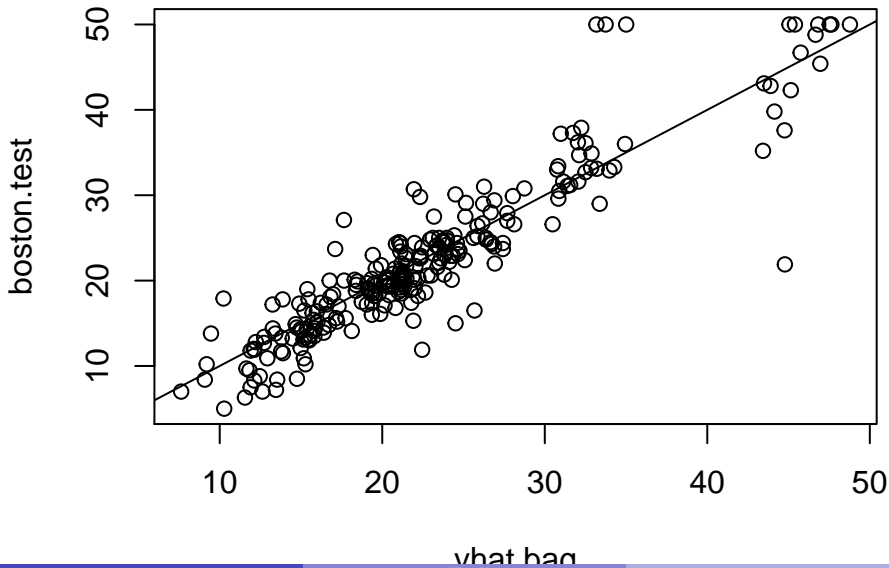
```
## No. of variables tried at each split: 13
```

```
##
```

```
##           Mean of squared residuals: 11.15723
```

```
##           % Var explained: 86.49
```

```
yhat.bag = predict(bag.boston,newdata=Boston[-train,])  
plot(yhat.bag, boston.test)  
abline(0,1)
```



```
bag.boston=randomForest(medv~.,data=Boston,subset=train,mtry=13,ntree=25)
yhat.bag = predict(bag.boston,newdata=Boston[-train,])
mean((yhat.bag-boston.test)^2)
```

```
## [1] 13.94835
```

```
set.seed(1)
rf.boston=randomForest(medv~.,data=Boston,subset=train,mtry=6,importance=TRUE)
yhat.rf = predict(rf.boston,newdata=Boston[-train,])
mean((yhat.rf-boston.test)^2)
```

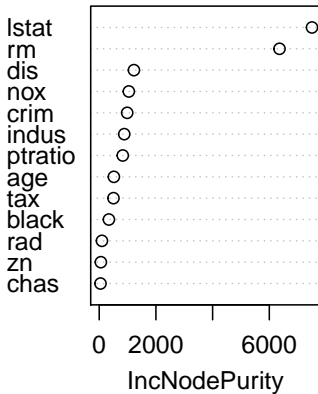
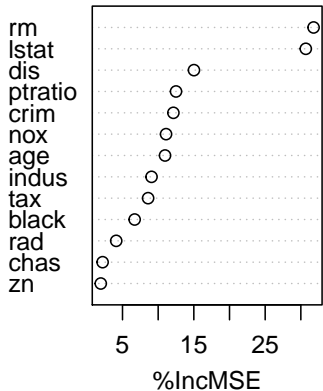
```
## [1] 11.66454
```

```
importance(rf.boston)
```

##		%IncMSE	IncNodePurity
##	crim	12.132320	986.50338
##	zn	1.955579	57.96945
##	indus	9.069302	882.78261
##	chas	2.210835	45.22941
##	nox	11.104823	1044.33776
##	rm	31.784033	6359.31971
##	age	10.962684	516.82969
##	dis	15.015236	1224.11605
##	rad	4.118011	95.94586
##	tax	8.587932	502.96719
##	prratio	12.503896	830.77523
##	black	6.702609	341.30361
##	lstat	30.695224	7505.73936

```
varImpPlot(rf.boston)
```

rf.boston



4. 부스팅

```
library(gbm)
```

```
## Loading required package: survival
```

```
## Loading required package: lattice
```

```
## Loading required package: splines
```

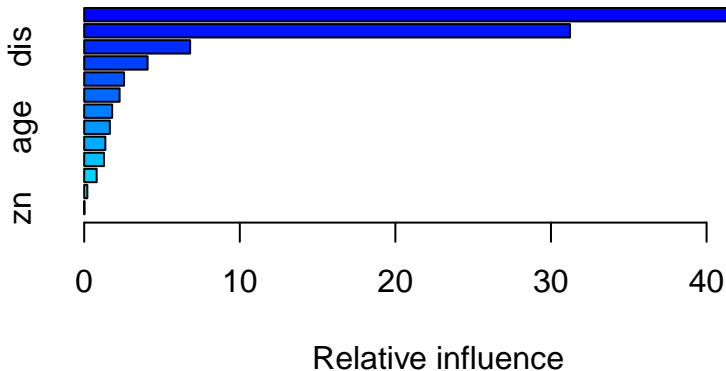
```
## Loading required package: parallel
```

```
## Loaded gbm 2.1.3
```

```
set.seed(1)
```

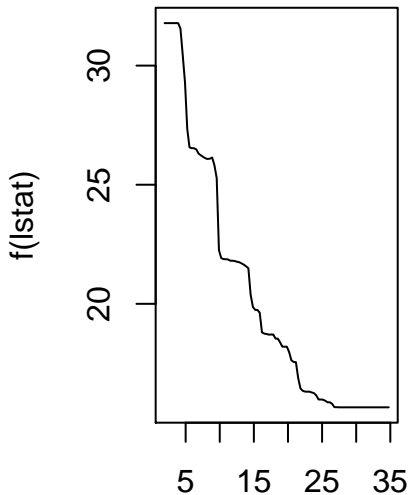
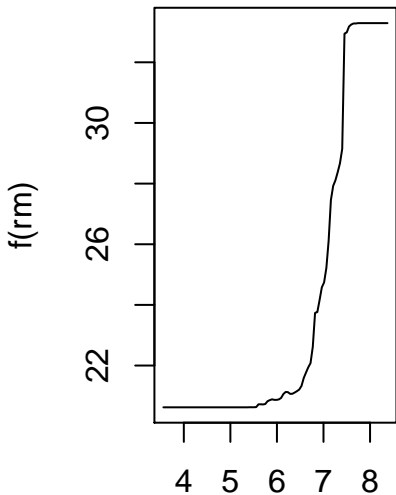
```
boost.boston=gbm(medv~.,data=Boston[train,],distribution="gaussian",n.trees=5000,int
```

```
summary(boost.boston)
```



```
##          var      rel.inf
## lstat    lstat 45.9627334
## rm       rm   31.2238187
## dis      dis   6.8087398
## crim     crim  4.0743784
```

```
par(mfrow=c(1,2))  
plot(boost.boston,i="rm")  
plot(boost.boston,i="lstat")
```




```
yhat.boost=predict(boost.boston,newdata=Boston[-train,],n.trees=5000)
mean((yhat.boost-boston.test)^2)
```

```
## [1] 11.84434
```

```
boost.boston=gbm(medv~.,data=Boston[train,],distribution="gaussian",n.trees=5000,int
yhat.boost=predict(boost.boston,newdata=Boston[-train,],n.trees=5000)
mean((yhat.boost-boston.test)^2)
```

```
## [1] 11.51109
```