# Lab 1: Subset Selection Methods

*Jonathan Bryan*

*June 1, 2018*

**6.5.1 Best Subset Selection**

```r
library(ISLR)
names(Hitters)
```

```
##  [1] "AtBat"     "Hits"      "HmRun"     "Runs"      "RBI"
##  [6] "Walks"     "Years"     "CAtBat"    "CHits"     "CHmRun"
## [11] "CRuns"     "CRBI"      "CWalks"    "League"    "Division"
## [16] "PutOuts"   "Assists"   "Errors"    "Salary"    "NewLeague"
```

```r
dim(Hitters)
```

```
## [1] 322  20
```

```r
sum(is.na(Hitters$Salary))
```

```
## [1] 59
```

```r
#omit observations with missing data
Hitters = na.omit(Hitters)
dim(Hitters)
```

```
## [1] 263  20
```

```r
sum(is.na(Hitters))
```

```
## [1] 0
```

```r
library(leaps)
```

```
## Warning: package 'leaps' was built under R version 3.4.4
```

```r
#best subset selection
regfit.full=regsubsets(Salary ~., Hitters)
summary(regfit.full)
```

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., Hitters)
## 19 Variables  (and intercept)
##           Forced in Forced out
## AtBat         FALSE      FALSE
## Hits          FALSE      FALSE
## HmRun         FALSE      FALSE
## Runs          FALSE      FALSE
## RBI           FALSE      FALSE
## Walks         FALSE      FALSE
## Years         FALSE      FALSE
## CAtBat        FALSE      FALSE
## CHits         FALSE      FALSE
## CHmRun        FALSE      FALSE
## CRuns         FALSE      FALSE
## CRBI          FALSE      FALSE
```

```
## CWalks           FALSE      FALSE
## LeagueN          FALSE      FALSE
## DivisionW        FALSE      FALSE
## PutOuts          FALSE      FALSE
## Assists          FALSE      FALSE
## Errors           FALSE      FALSE
## NewLeagueN       FALSE      FALSE
## 1 subsets of each size up to 8
## Selection Algorithm: exhaustive
##          AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns
## 1  ( 1 ) " "   " "  " "   " "  " " " "   " "   " "    " "   " "    " "
## 2  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 3  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 4  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 5  ( 1 ) "*"   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 6  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    " "
## 7  ( 1 ) " "   "*"  " "   " "  " " "*"   " "   "*"    "*"   "*"    " "
## 8  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   "*"    "*"
##          CRBI CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
## 1  ( 1 ) "*"  " "    " "     " "       " "     " "     " "    " "
## 2  ( 1 ) "*"  " "    " "     " "       " "     " "     " "    " "
## 3  ( 1 ) "*"  " "    " "     " "       "*"     " "     " "    " "
## 4  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 5  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 6  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 7  ( 1 ) " "  " "    " "     "*"       "*"     " "     " "    " "
## 8  ( 1 ) " "  "*"    " "     "*"       "*"     " "     " "    " "
```

```r
regfit.full = regsubsets(Salary ~., data=Hitters, nvmax=19)
reg.summary = summary(regfit.full)
names(reg.summary)
```

```
## [1] "which"  "rsq"    "rss"    "adjr2"  "cp"     "bic"    "outmat" "obj"
```

```r
#analysis of model selection criteria
round(reg.summary$rsq,3)
```

```
##  [1] 0.321 0.425 0.451 0.475 0.491 0.509 0.514 0.529 0.535 0.540 0.543
## [12] 0.544 0.544 0.545 0.545 0.546 0.546 0.546 0.546
```

```r
par(mfrow=c(2,2))
plot(reg.summary$rss,xlab="Number of Variables", ylab="RSS", type = "l")
plot(reg.summary$adjr2,xlab="Number of Variables", ylab="Adjusted RSq", type = "l")
which.max (reg.summary$adjr2)
```
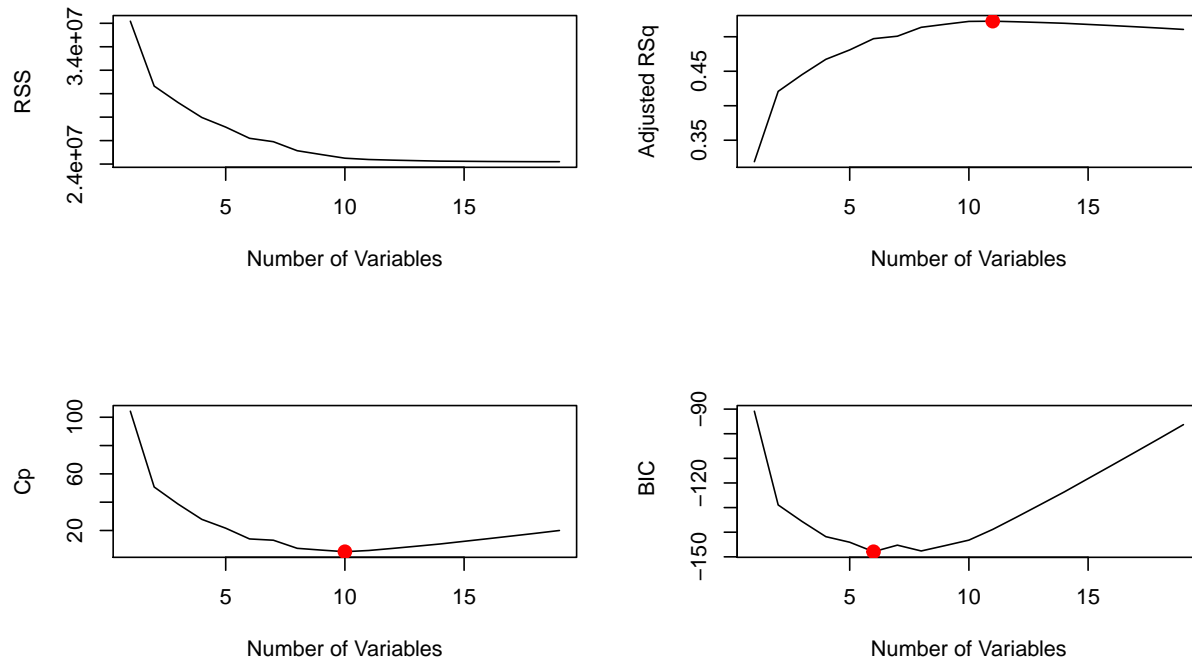
```
## [1] 11
```

```r
points (11, reg.summary$adjr2[11], col ="red",cex =2, pch =20)
plot(reg.summary$cp ,xlab =" Number of Variables ",ylab="Cp",type="l")
which.min (reg.summary$cp)
```
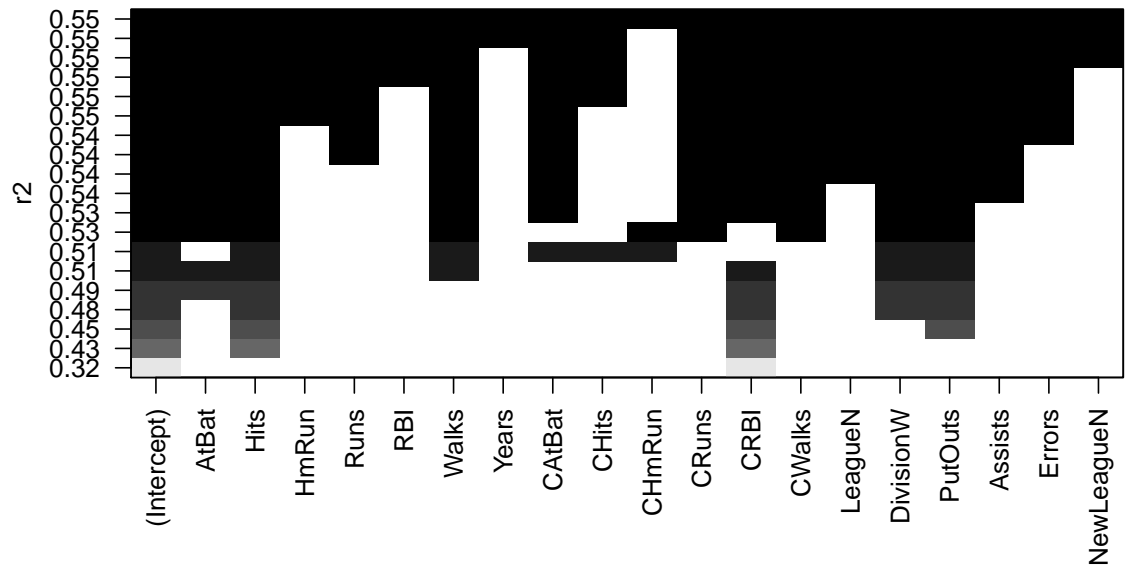
```
## [1] 10
```

```r
points(10, reg.summary$cp [10], col ="red",cex =2, pch =20)
which.min(reg.summary$bic)
```
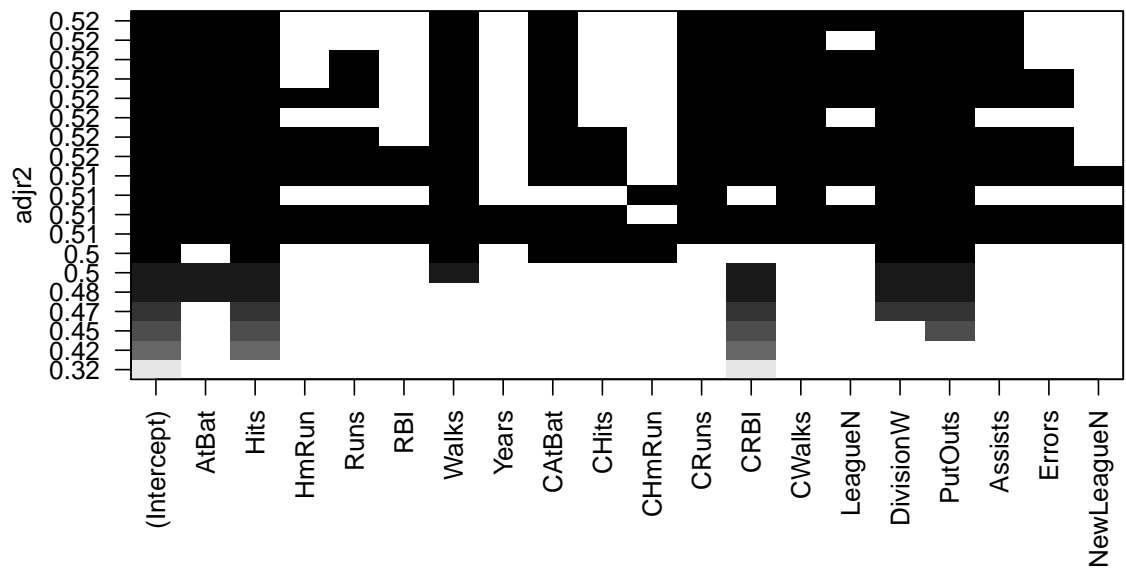
```
## [1] 6
```

```r
plot(reg.summary$bic ,xlab=" Number of Variables ",ylab=" BIC", type="l")
points (6, reg.summary$bic [6], col =" red",cex =2, pch =20)
```
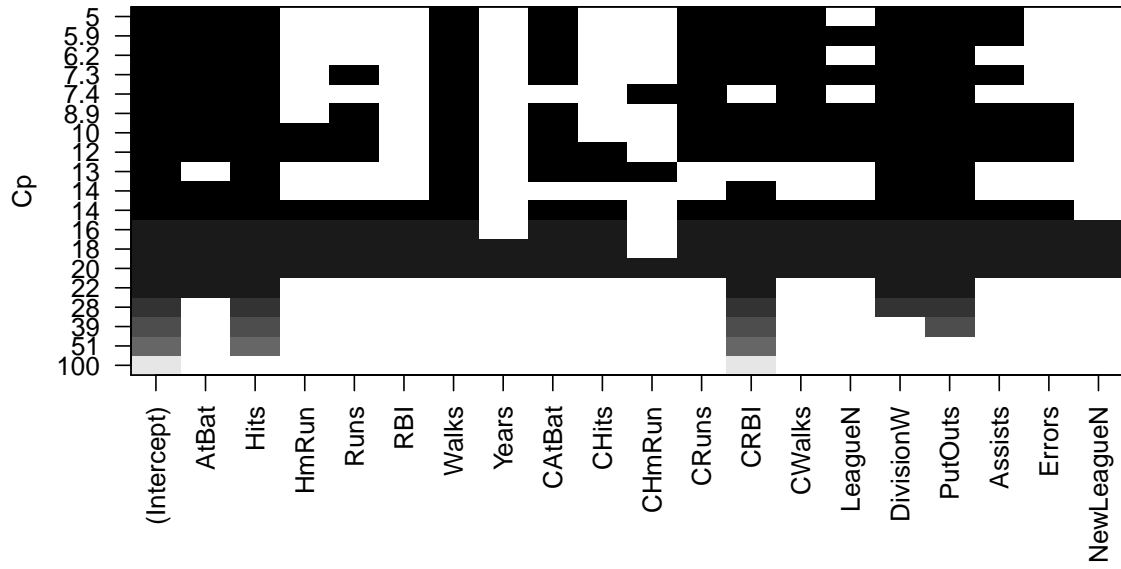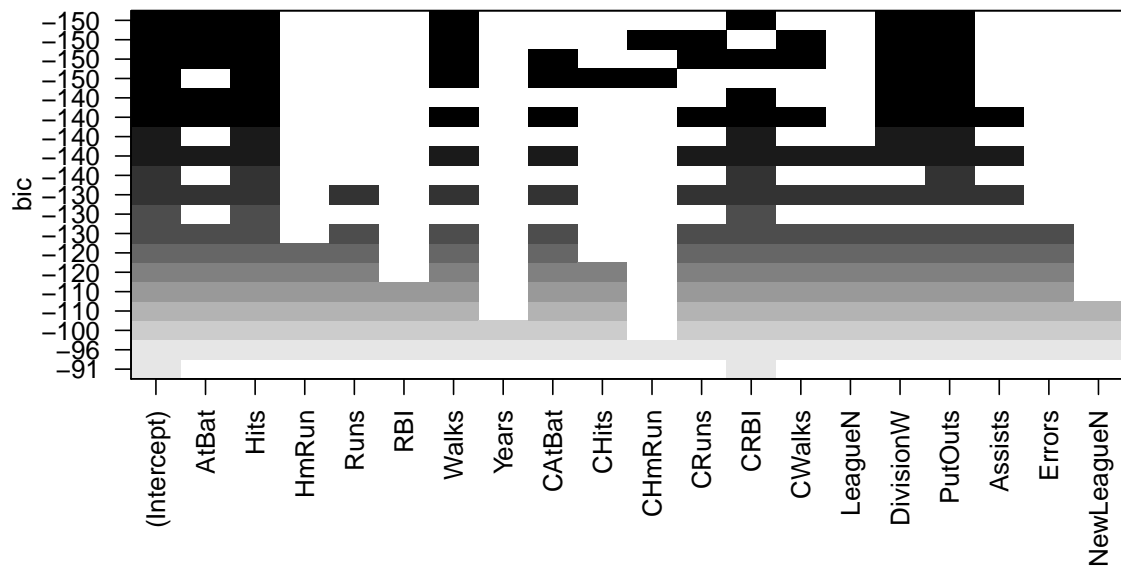


```r
plot(regfit.full,scale="r2")
```

```r
plot(regfit.full ,scale ="adjr2")
```

```r
plot(regfit.full ,scale ="Cp")
```



```r
plot(regfit.full ,scale ="bic")
```

```r
#coefficients of six variable model
coef(regfit.full,6)
```

```
##   (Intercept)         AtBat          Hits         Walks          CRBI
##    91.5117981    -1.8685892     7.6043976     3.6976468     0.6430169
##      DivisionW       PutOuts
##  -122.9515338     0.2643076
```

**Forward and Backward Stepwise Selection**

```r
regfit.fwd = regsubsets(Salary ~ ., data = Hitters, method="forward", nvmax = 19)
summary(regfit.fwd)
```

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., data = Hitters, method = "forward",
##     nvmax = 19)
## 19 Variables  (and intercept)
##            Forced in Forced out
## AtBat          FALSE      FALSE
## Hits           FALSE      FALSE
## HmRun          FALSE      FALSE
## Runs           FALSE      FALSE
## RBI            FALSE      FALSE
## Walks          FALSE      FALSE
## Years          FALSE      FALSE
## CAtBat         FALSE      FALSE
## CHits          FALSE      FALSE
## CHmRun         FALSE      FALSE
## CRuns          FALSE      FALSE
## CRBI           FALSE      FALSE
## CWalks         FALSE      FALSE
## LeagueN        FALSE      FALSE
## DivisionW      FALSE      FALSE
## PutOuts        FALSE      FALSE
## Assists        FALSE      FALSE
## Errors         FALSE      FALSE
## NewLeagueN     FALSE      FALSE
## 1 subsets of each size up to 19
## Selection Algorithm: forward
##           AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns
## 1  ( 1 )  " "   " "  " "   " "  " " " "   " "   " "    " "   " "    " "
## 2  ( 1 )  " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 3  ( 1 )  " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 4  ( 1 )  " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 5  ( 1 )  "*"   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "
## 6  ( 1 )  "*"   "*"  " "   " "  " " " "   "*"   " "    " "   " "    " "
## 7  ( 1 )  "*"   "*"  " "   " "  " " " "   "*"   " "    " "   " "    " "
## 8  ( 1 )  "*"   "*"  " "   " "  " " " "   "*"   " "    " "   " "    "*"
## 9  ( 1 )  "*"   "*"  " "   " "  " " " "   "*"   " "    "*"   " "    "*"
## 10  ( 1 ) "*"   "*"  " "   " "  " " " "   "*"   " "    "*"   " "    "*"
## 11  ( 1 ) "*"   "*"  " "   " "  " " " "   "*"   " "    "*"   " "    "*"
## 12  ( 1 ) "*"   "*"  " "   "*"  " " " "   "*"   " "    "*"   " "    "*"
## 13  ( 1 ) "*"   "*"  " "   "*"  " " " "   "*"   " "    "*"   " "    "*"
```

```
## 14  ( 1 ) "*"   "*"   "*"   "*"   " " "*"   " "    "*"     " "    " "    "*"
## 15  ( 1 ) "*"   "*"   "*"   "*"   " " "*"   " "    "*"     "*"    " "    "*"
## 16  ( 1 ) "*"   "*"   "*"   "*"   "*" "*"   " "    "*"     "*"    " "    "*"
## 17  ( 1 ) "*"   "*"   "*"   "*"   "*" "*"   " "    "*"     "*"    " "    "*"
## 18  ( 1 ) "*"   "*"   "*"   "*"   "*" "*"   "*"    "*"     "*"    " "    "*"
## 19  ( 1 ) "*"   "*"   "*"   "*"   "*" "*"   "*"    "*"     "*"    "*"    "*"
##           CRBI CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
## 1  ( 1 ) "*"  " "    " "     " "       " "     " "     " "    " "
## 2  ( 1 ) "*"  " "    " "     " "       " "     " "     " "    " "
## 3  ( 1 ) "*"  " "    " "     " "       "*"     " "     " "    " "
## 4  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 5  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 6  ( 1 ) "*"  " "    " "     "*"       "*"     " "     " "    " "
## 7  ( 1 ) "*"  "*"    " "     "*"       "*"     " "     " "    " "
## 8  ( 1 ) "*"  "*"    " "     "*"       "*"     " "     " "    " "
## 9  ( 1 ) "*"  "*"    " "     "*"       "*"     " "     " "    " "
## 10 ( 1 ) "*"  "*"    " "     "*"       "*"     "*"     " "    " "
## 11 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     " "    " "
## 12 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     " "    " "
## 13 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 14 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 15 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 16 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 17 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
## 18 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
## 19 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
```

```r
regfit.bwd = regsubsets(Salary ~ ., data = Hitters, method = "backward", nvmax = 19)
summary(regfit.bwd)
```

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., data = Hitters, method = "backward",
##     nvmax = 19)
## 19 Variables  (and intercept)
##            Forced in Forced out
## AtBat          FALSE      FALSE
## Hits           FALSE      FALSE
## HmRun          FALSE      FALSE
## Runs           FALSE      FALSE
## RBI            FALSE      FALSE
## Walks          FALSE      FALSE
## Years          FALSE      FALSE
## CAtBat         FALSE      FALSE
## CHits          FALSE      FALSE
## CHmRun         FALSE      FALSE
## CRuns          FALSE      FALSE
## CRBI           FALSE      FALSE
## CWalks         FALSE      FALSE
## LeagueN        FALSE      FALSE
## DivisionW      FALSE      FALSE
## PutOuts        FALSE      FALSE
## Assists        FALSE      FALSE
## Errors         FALSE      FALSE
## NewLeagueN     FALSE      FALSE
## 1 subsets of each size up to 19
```

```
## Selection Algorithm: backward
##          AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns
## 1  ( 1 ) " "   " "  " "   " "  " " " "   " "   " "    " "   " "    "*"
## 2  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    "*"
## 3  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    "*"
## 4  ( 1 ) "*"   "*"  " "   " "  " " " "   " "   " "    " "   " "    "*"
## 5  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    "*"
## 6  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    "*"
## 7  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    "*"
## 8  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    "*"
## 9  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   "*"    " "   " "    "*"
## 10 ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   "*"    " "   " "    "*"
## 11 ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   "*"    " "   " "    "*"
## 12 ( 1 ) "*"   "*"  " "   "*"  " " "*"   " "   "*"    " "   " "    "*"
## 13 ( 1 ) "*"   "*"  " "   "*"  " " "*"   " "   "*"    " "   " "    "*"
## 14 ( 1 ) "*"   "*"  "*"   "*"  " " "*"   " "   "*"    " "   " "    "*"
## 15 ( 1 ) "*"   "*"  "*"   "*"  " " "*"   " "   "*"    "*"   " "    "*"
## 16 ( 1 ) "*"   "*"  "*"   "*"  "*" "*"   " "   "*"    "*"   " "    "*"
## 17 ( 1 ) "*"   "*"  "*"   "*"  "*" "*"   " "   "*"    "*"   " "    "*"
## 18 ( 1 ) "*"   "*"  "*"   "*"  "*" "*"   "*"   "*"    "*"   " "    "*"
## 19 ( 1 ) "*"   "*"  "*"   "*"  "*" "*"   "*"   "*"    "*"   "*"    "*"
##          CRBI CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
## 1  ( 1 ) " "  " "    " "     " "       " "     " "     " "    " "
## 2  ( 1 ) " "  " "    " "     " "       " "     " "     " "    " "
## 3  ( 1 ) " "  " "    " "     " "       "*"     " "     " "    " "
## 4  ( 1 ) " "  " "    " "     " "       "*"     " "     " "    " "
## 5  ( 1 ) " "  " "    " "     " "       "*"     " "     " "    " "
## 6  ( 1 ) " "  " "    " "     "*"       "*"     " "     " "    " "
## 7  ( 1 ) " "  "*"    " "     "*"       "*"     " "     " "    " "
## 8  ( 1 ) "*"  "*"    " "     "*"       "*"     " "     " "    " "
## 9  ( 1 ) "*"  "*"    " "     "*"       "*"     " "     " "    " "
## 10 ( 1 ) "*"  "*"    " "     "*"       "*"     "*"     " "    " "
## 11 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     " "    " "
## 12 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     " "    " "
## 13 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 14 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 15 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 16 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    " "
## 17 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
## 18 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
## 19 ( 1 ) "*"  "*"    "*"     "*"       "*"     "*"     "*"    "*"
```

```r
round(coef(regfit.full,7),3)
```

```
## (Intercept)        Hits        Walks      CAtBat        CHits      CHmRun
##      79.451       1.283        3.227      -0.375        1.496       1.442
##    DivisionW     PutOuts
##    -129.987       0.237
```

```r
round(coef(regfit.fwd,7),3)
```

```
## (Intercept)        AtBat         Hits       Walks         CRBI      CWalks
##     109.787       -1.959        7.450       4.913        0.854      -0.305
##    DivisionW     PutOuts
##    -127.122       0.253
```

```r
round(coef(regfit.bwd,7),3)
```

```
## (Intercept)        AtBat         Hits        Walks        CRuns       CWalks
##     105.649       -1.976        6.757        6.056        1.129       -0.716
##   DivisionW      PutOuts
##    -116.169        0.303
```

```r
set.seed(1)
train = sample(c(TRUE,FALSE), size = nrow(Hitters), rep = TRUE)
test = (!train)

#fit training model
regfit.best = regsubsets(Salary ~ ., data = Hitters, subset = train, nvmax=19)

#validation set
test.mat = model.matrix(Salary ~ ., data = Hitters[test,])
val.errors = rep(NA,19)
for(i in 1:19){
  coefi = coef(regfit.best, id=i)
  pred = test.mat[,names(coefi)]%*%coefi
  val.errors[i] = mean((Hitters$Salary[test]-pred)^2)
}
val.errors
```

```
##  [1] 220968.0 169157.1 178518.2 163426.1 168418.1 171270.6 162377.1
##  [8] 157909.3 154055.7 148162.1 151156.4 151742.5 152214.5 157358.7
## [15] 158541.4 158743.3 159972.7 159859.8 160105.6
```

```r
#observe that the 10 variable model has the lowest test MSE
which.min(val.errors)
```

```
## [1] 10
```

```r
coef(regfit.best, 10)
```

```
## (Intercept)        AtBat         Hits        Walks       CAtBat        CHits
## -80.2751499   -1.4683816    7.1625314    3.6430345   -0.1855698    1.1053238
##      CHmRun       CWalks      LeagueN     DivisionW      PutOuts
##   1.3844863   -0.7483170   84.5576103  -53.0289658    0.2381662
```

```r
#prediction function for subsets
predict.regsubsets = function(object,newdata,id,...){
  form = as.formula(object$call[[2]])
  mat = model.matrix(form, newdata)
  coefi = coef(object,id = id)
  xvars = names(coefi)
  mat[,xvars]%*%coefi
}

#fit the best subset 10 variable model and observe coefficients
regfit.best = regsubsets(Salary ~., data = Hitters, nvmax = 19)
coef(regfit.best, 10)
```

```
## (Intercept)        AtBat         Hits        Walks       CAtBat
## 162.5354420   -2.1686501    6.9180175    5.7732246   -0.1300798
##        CRuns         CRBI       CWalks     DivisionW      PutOuts
##    1.4082490    0.7743122   -0.8308264  -112.3800575    0.2973726
```

```
##    Assists
##   0.2831680
```
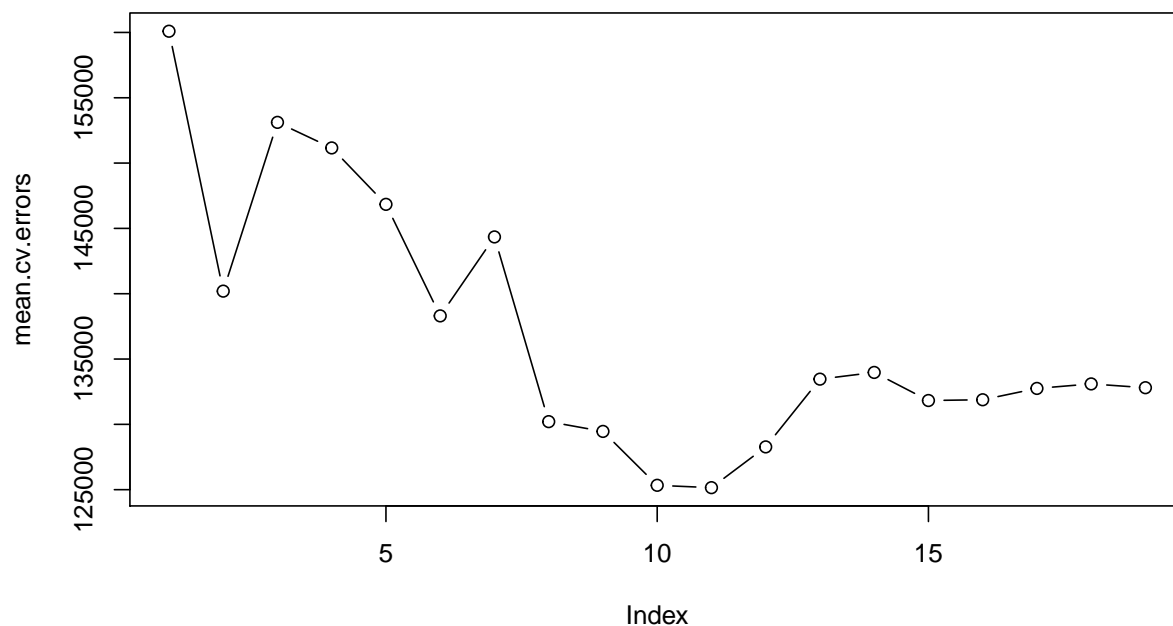
```r
#k-fold CV for best subset
k=10
set.seed(1)
folds = sample(1:k, nrow(Hitters), replace = TRUE)
cv.errors = matrix(NA,k,19,dimnames = list(NULL, paste(1:19)))

#for each k-fold, run best subset for 1 to 19 variables
for(j in 1:k){
  best.fit = regsubsets(Salary ~., data = Hitters[folds!=j,], nvmax =19)
  for (i in 1:19){
    pred = predict.regsubsets(best.fit, Hitters[folds == j,], id=i)
    cv.errors[j,i] = mean((Hitters$Salary[folds ==j] - pred)^2)
  }
}
```

```r
mean.cv.errors = apply(cv.errors,2,mean)
mean.cv.errors
```

```
##        1        2        3        4        5        6        7        8
## 160093.5 140196.8 153117.0 151159.3 146841.3 138302.6 144346.2 130207.7
##        9       10       11       12       13       14       15       16
## 129459.6 125334.7 125153.8 128273.5 133461.0 133974.6 131825.7 131882.8
##       17       18       19
## 132750.9 133096.2 132804.7
```

```r
par(mfrow = c(1,1))
plot(mean.cv.errors,type="b")
```

```
reg.best = regsubsets(Salary ~., data = Hitters, nvmax =19)
coef(reg.best,11)
```

```
## (Intercept)        AtBat          Hits         Walks        CAtBat
## 135.7512195   -2.1277482     6.9236994     5.6202755    -0.1389914
##        CRuns         CRBI        CWalks       LeagueN       DivisionW
##    1.4553310    0.7852528    -0.8228559    43.1116152  -111.1460252
##      PutOuts       Assists
##    0.2894087    0.2688277
```