# CH6_Regression_ex1

*Philip oh*

```
library(ISLR)
library(leaps)
```

```
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"
```

```
dim(Hitters)
```

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

```
str(Hitters)
```

```
## 'data.frame':    322 obs. of  20 variables:
##  $ AtBat    : int  293 315 479 496 321 594 185 298 323 401 ...
##  $ Hits     : int  66 81 130 141 87 169 37 73 81 92 ...
##  $ HmRun    : int  1 7 18 20 10 4 1 0 6 17 ...
##  $ Runs     : int  30 24 66 65 39 74 23 24 26 49 ...
##  $ RBI      : int  29 38 72 78 42 51 8 24 32 66 ...
##  $ Walks    : int  14 39 76 37 30 35 21 7 8 65 ...
##  $ Years    : int  1 14 3 11 2 11 2 3 2 13 ...
##  $ CAtBat   : int  293 3449 1624 5628 396 4408 214 509 341 5206 ...
##  $ CHits    : int  66 835 457 1575 101 1133 42 108 86 1332 ...
##  $ CHmRun   : int  1 69 63 225 12 19 1 0 6 253 ...
##  $ CRuns    : int  30 321 224 828 48 501 30 41 32 784 ...
##  $ CRBI     : int  29 414 266 838 46 336 9 37 34 890 ...
##  $ CWalks   : int  14 375 263 354 33 194 24 12 8 866 ...
##  $ League   : Factor w/ 2 levels "A","N": 1 2 1 2 2 1 2 1 2 1 ...
##  $ Division : Factor w/ 2 levels "E","W": 1 2 2 1 1 2 1 2 2 1 ...
##  $ PutOuts  : int  446 632 880 200 805 282 76 121 143 0 ...
##  $ Assists  : int  33 43 82 11 40 421 127 283 290 0 ...
##  $ Errors   : int  20 10 14 3 4 25 7 9 19 0 ...
##  $ Salary   : num  NA 475 480 500 91.5 750 70 100 75 1100 ...
##  $ NewLeague: Factor w/ 2 levels "A","N": 1 2 1 2 2 1 1 1 2 1 ...
```

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

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

- 자료를 살펴보니 Salary에 NA가 있다.

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

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

```
Hitters = na.omit(Hitters)
sum(is.na(Hitters))
```

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

- NA를 모두 제거했다.

```
regfit.full = regsubsets(Salary ~ ., data = Hitters)
summary(regfit.full)
```

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., data = 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 ) " "  "*"    " "     "*"       "*"     " "     " "    " "
```

- `regsubsets` 는 변수의 개수에 따른 최적의 모형을 반환한다.
- `regsubsets` 의 옵션 중 `force.in` 과 `force.out` 은 반드시 모형에 들어가야하는 혹은 빠져야 하는 변수들의 인덱스를 지정한다. `summary` 에서 모형의 크기가 8개까지인 것만 보여주는데 이것을 바꾸려면 `nvmax` 옵션을 쓰면 된다.

```
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"
```

- 19개의 변수까지 최적의 모형을 반환한다. `reg.summary$rsq` 는 변수의 개수에 따른 최적의 모형의 R Square값을 갖고 있다.

```
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
```

```
points(11, reg.summary$adjr2[11], col = "red", cex = 2, pch = 20)

which.min(reg.summary$cp)
```
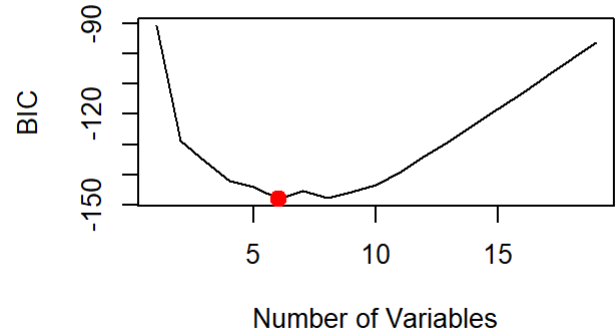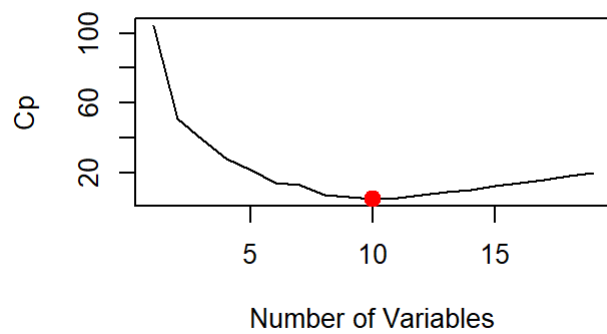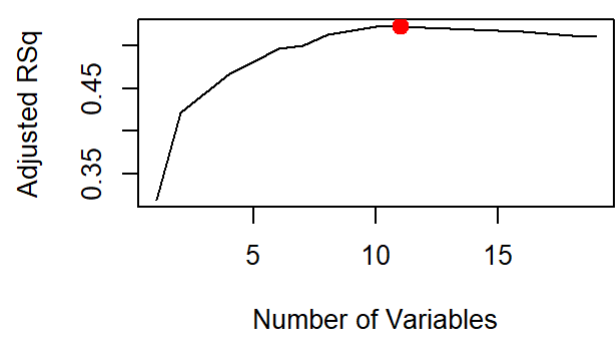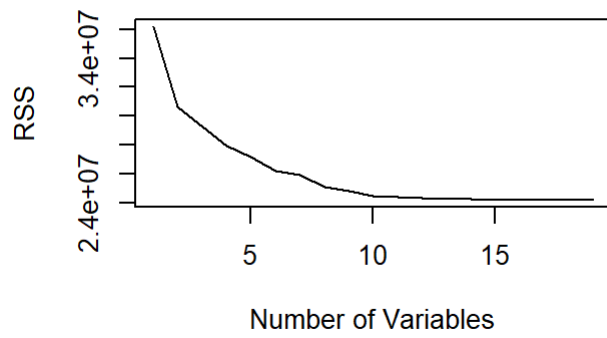
```
## [1] 10
```

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

which.min(reg.summary$bic)
```

```
## [1] 6
```

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

- max 는 최대값 그 자체를 반환하고, which.max 는 최대값의 위치를 반환한다.
- 변수의 개수에 따른 최적 모형의 RSS, adjusted R2, Cp, BIC 값을 그림으로 그리고 최적의 모형을 표시했다.

```
par(mfrow=c(1, 1))
plot(regfit.full, scale="r2")
```

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

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



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

```
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
```
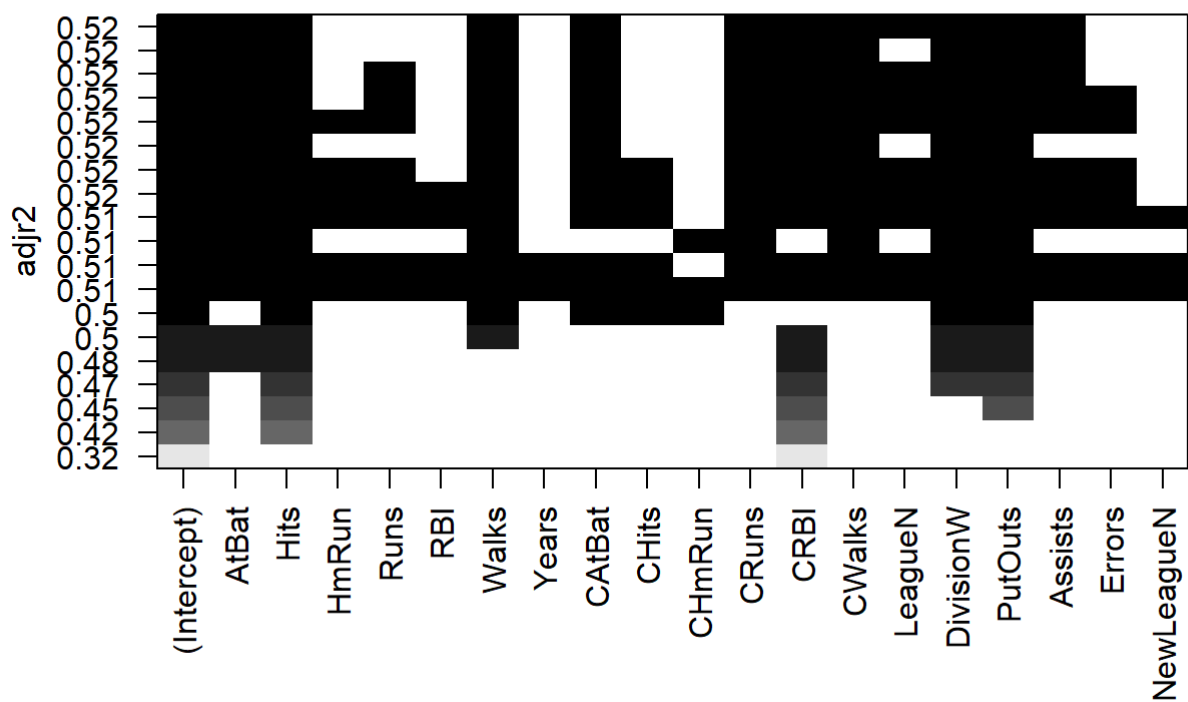
- 이 그림들은 각 기준에 따라 들어가는 변수들을 검은색 박스로 표현했다. 선택되지 않은 변수는 흰색박스로 표시된다. bic의 기준으로 최적의 모형은 6개의 변수(절편을 포함하면 7개)를 포함하는 모형이다. 이 모형의 회귀계수를 알아보기 위해 마지막 명령어를 썼다.

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

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., data = Hitters, nvmax = 19, method = "forward")
## 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 ) "*"  "*"   "*"    "*"      "*"     "*"     "*"      "*"
```

- 전진선택법을 이용할 때 선택되는 변수들을 *로 표시한다.

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

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., data = Hitters, nvmax = 19, method = "backward")
## 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 ) "*" "*"  "*"   "*"    "*"   "*"   "*"   "*"
```

- 후진선택법으로 변수를 선택할 때의 결과를 보여준다.

```
coef(regfit.full, 7)
```

```
## (Intercept)        Hits        Walks       CAtBat        CHits
##  79.4509472   1.2833513    3.2274264   -0.3752350    1.4957073
##      CHmRun   DivisionW      PutOuts
##   1.4420538 -129.9866432    0.2366813
```

```
coef(regfit.fwd, 7)
```

```
## (Intercept)        AtBat        Hits        Walks        CRBI
## 109.7873062   -1.9588851   7.4498772    4.9131401    0.8537622
##      CWalks   DivisionW      PutOuts
##   -0.3053070 -127.1223928    0.2533404
```

```
coef(regfit.bwd, 7)
```

```
## (Intercept)        AtBat        Hits        Walks        CRuns
## 105.6487488   -1.9762838   6.7574914    6.0558691    1.1293095
##      CWalks   DivisionW      PutOuts
##   -0.7163346 -116.1692169    0.3028847
```

- all subset selection, 전진선택법, 후진선택법을 이용할 때 변수 7개의 최적 모형이 다 다르다. 변수 6개까지의 모형은 세 가지 방법이 모두 같다.