# Hitters Notebook

The hitters dataset consists of 322 observations of 21 variables with the following information - X (name), At-Bat, Hits, HmRun (home runs), Runs, RBI, Walks, Years, CAtBat, CHits, CHmRun, CRuns, CRBI, CWalks, League, Division, PutOuts, Assists, Errors, Salary, New League. Here League, Division and NewLeagues are factor variabes with 2 categories. We drop rows with missing entries and are left with 263 observations.

```r
rm(list=ls())
hitter <- read.csv("hitters.csv")
# str(hitter)
# summary(hitter)
hitter <- na.omit(hitter)
# str(hitter)
# summary(hitter)
```

The leaps package in R does subset selection with the `regsubsets` function. By default, the maximum number of subsets, this function uses is 8 (ie. number of variables in the model, till M8). We extend this to do a complete subset selection by changing the default value of `nvmax` argument in this function. Note that CRBI is in the model with 1 to 6 variables but not in the model with 7 and 8 variables.

```r
#install.packages("leaps")
library(leaps)
?regsubsets
# select the columns comprising of the variables
hitters <- hitter[, 2:21]
# predict salary with all the variables
# default: best subset selection algorithm
# method = "exhaustive"
model1 <- regsubsets(Salary ~ ., hitters)
summary(model1)
```

```
## 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 CRBI
## 1  ( 1 ) " "   " "  " "   " "  " " " "   " "   " "    " "   " "    " "   "*"
## 2  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "   "*"
## 3  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "   "*"
## 4  ( 1 ) " "   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "   "*"
## 5  ( 1 ) "*"   "*"  " "   " "  " " " "   " "   " "    " "   " "    " "   "*"
## 6  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   " "    " "   "*"
## 7  ( 1 ) " "   "*"  " "   " "  " " "*"   " "   "*"    "*"   "*"    " "   " "
## 8  ( 1 ) "*"   "*"  " "   " "  " " "*"   " "   " "    " "   "*"    "*"   " "
##          CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
## 1  ( 1 ) " "    " "     " "       " "     " "     " "    " "
## 2  ( 1 ) " "    " "     " "       " "     " "     " "    " "
## 3  ( 1 ) " "    " "     " "       "*"     " "     " "    " "
## 4  ( 1 ) " "    " "     "*"       "*"     " "     " "    " "
## 5  ( 1 ) " "    " "     "*"       "*"     " "     " "    " "
## 6  ( 1 ) " "    " "     "*"       "*"     " "     " "    " "
## 7  ( 1 ) " "    " "     "*"       "*"     " "     " "    " "
## 8  ( 1 ) "*"    " "     "*"       "*"     " "     " "    " "
```

```r
# FORCED IN & FORCED OUT: to specify which variable i must include/exclude in the model
# by default, none is specified so all values return FALSE
# 1. choosing among all of the 19 variables available, and select the one variable with the minimal MSE
# CRBI got kicked out in the 7th selection

model2 <- regsubsets(Salary ~ ., hitters, nvmax = 19)
summary(model2)
```

```
## Subset selection object
## Call: regsubsets.formula(Salary ~ ., hitters, 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
```
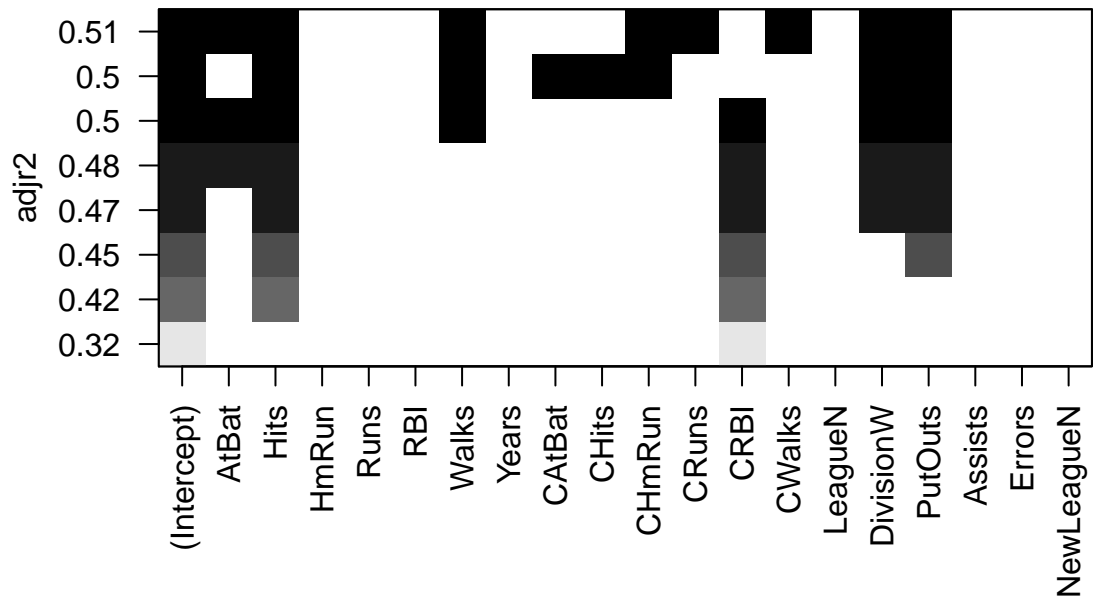
```
## DivisionW      FALSE       FALSE
## PutOuts         FALSE       FALSE
## Assists         FALSE       FALSE
## Errors          FALSE       FALSE
## NewLeagueN      FALSE       FALSE
## 1 subsets of each size up to 19
## Selection Algorithm: exhaustive
##           AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns CRBI
## 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 )  "*"   "*"  "*"   "*"  "*" "*"   "*"   "*"    "*"   "*"    "*"   "*"
##           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 )  "*"    "*"     "*"       "*"     "*"     "*"    "*"
```
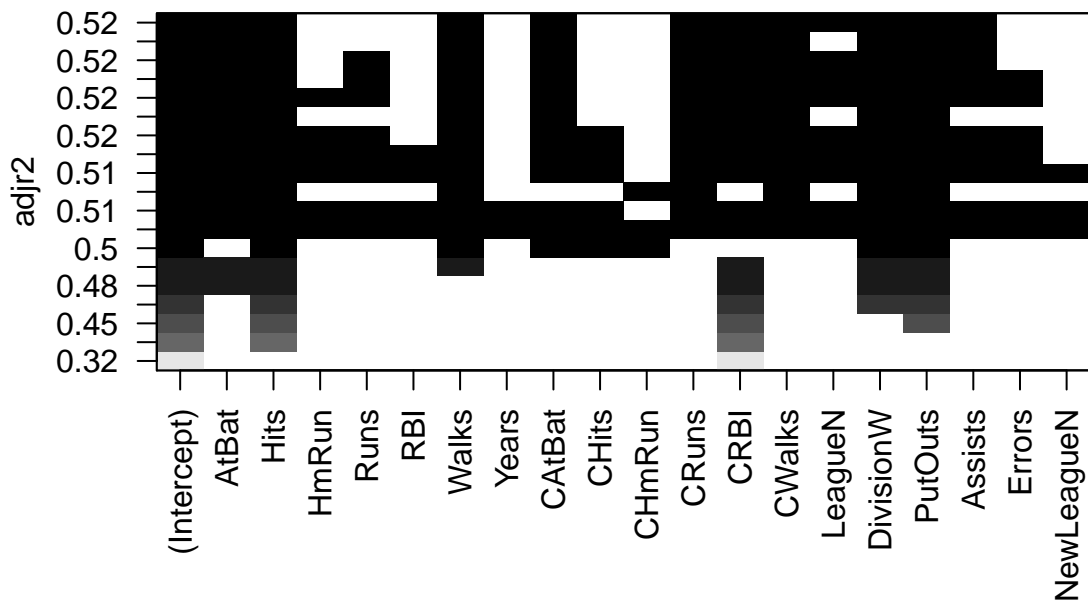
```r
plot(model1, scale = "adjr2")
```

```
plot(model2, scale = "adjr2")
```

```r
# gives what is the influence on the model at different levels of r2

# values that i can get from running summary(model2)
names(summary(model2))
```
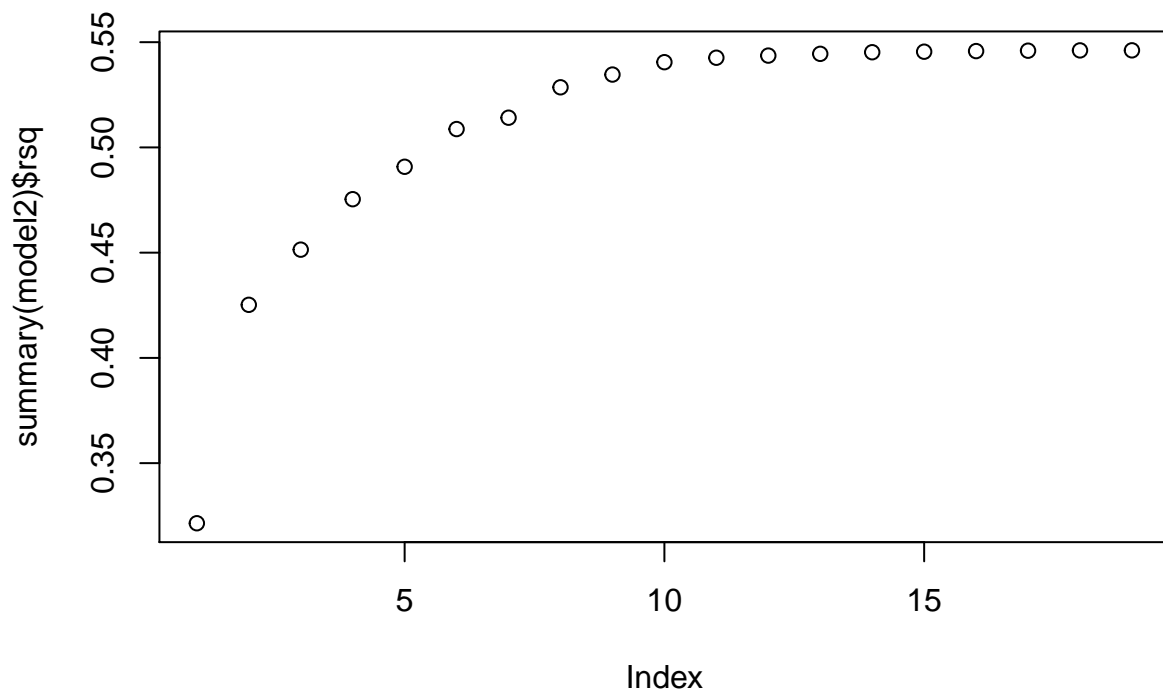
```
## [1] "which"  "rsq"    "rss"    "adjr2"  "cp"    "bic"    "outmat" "obj"
```
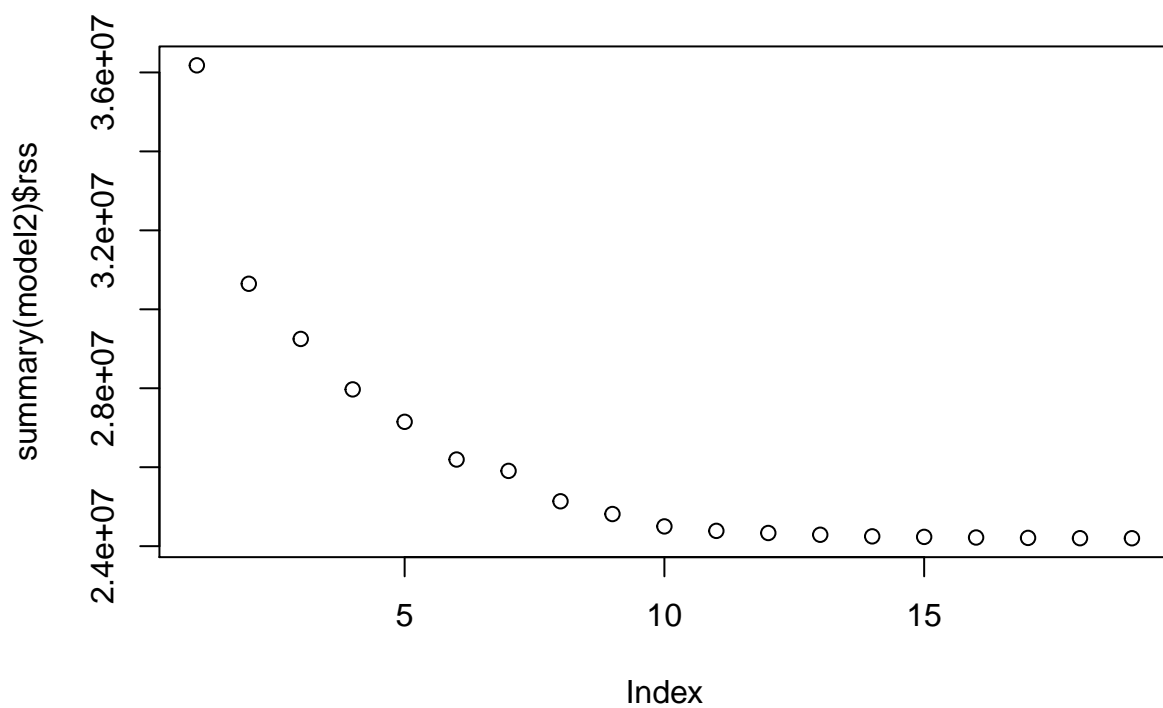
```r
# returns values of the best rsquare based on the variable size of subset selection, in ascending order
summary(model2)$rsq
```

```
##  [1] 0.3214501 0.4252237 0.4514294 0.4754067 0.4908036 0.5087146 0.5141227
##  [8] 0.5285569 0.5346124 0.5404950 0.5426153 0.5436302 0.5444570 0.5452164
## [15] 0.5454692 0.5457656 0.5459518 0.5460945 0.5461159
```
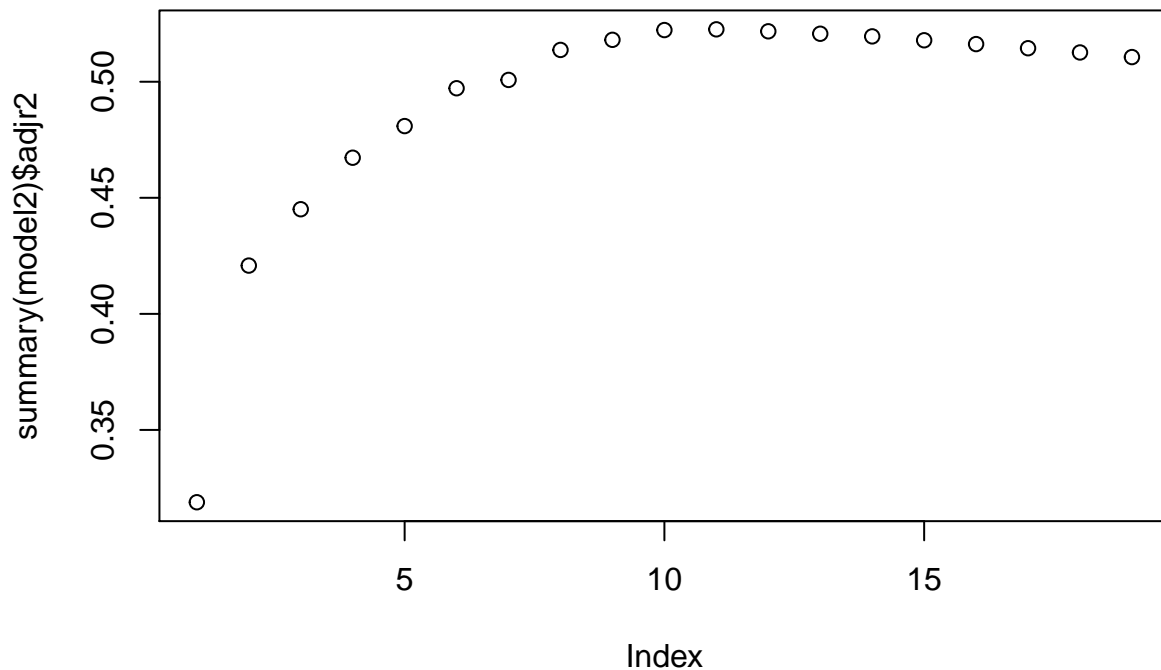
```r
plot(summary(model2)$rsq)
```

```
# value of r2 increases as the number of predictors increases
# hence we tend to not use r2 value to choose model

# rss: residual sum of squares
plot(summary(model2)$rss)
```

```
# value decreases as number of predictors increases as well

# adjr2: adjusted e2
plot(summary(model2)$adjr2)
```

```
which.max(summary(model2)$adjr2)
```

```
## [1] 11
```

```
# obtain the coefficients based of the best model (ie. lowest MSE)
coef(model2, 11)
```

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

The figures indicate that R-squared increase as the number of variables in the subset increases and likewise the residual sum of squared (sum of squared errors) decreases as the size of the subsets increases. On the other hand the adjusted R-squared increases first and then decreases.

Forward stepwise selection: In this example, the best model identified by the forward stepwise selection is the same as that obtained by the best subset selection. It is also possible to run this algorithm using a backward method where you drop variables one a time rather add. In general, the solutions from these two methods can be different.

```
# forward selection
model3 <- regsubsets(Salary ~ ., data = hitters, nvmax = 19, method = "forward")
which.max(summary(model3)$adjr2)
```

```
## [1] 11
```

```
coef(model3, 11)
```

```
## (Intercept)        AtBat          Hits         Walks        CAtBat         CRuns
## 135.7512195   -2.1277482     6.9236994     5.6202755    -0.1389914     1.4553310
##        CRBI       CWalks       LeagueN      DivisionW       PutOuts        Assists
##   0.7852528   -0.8228559    43.1116152  -111.1460252     0.2894087     0.2688277
```
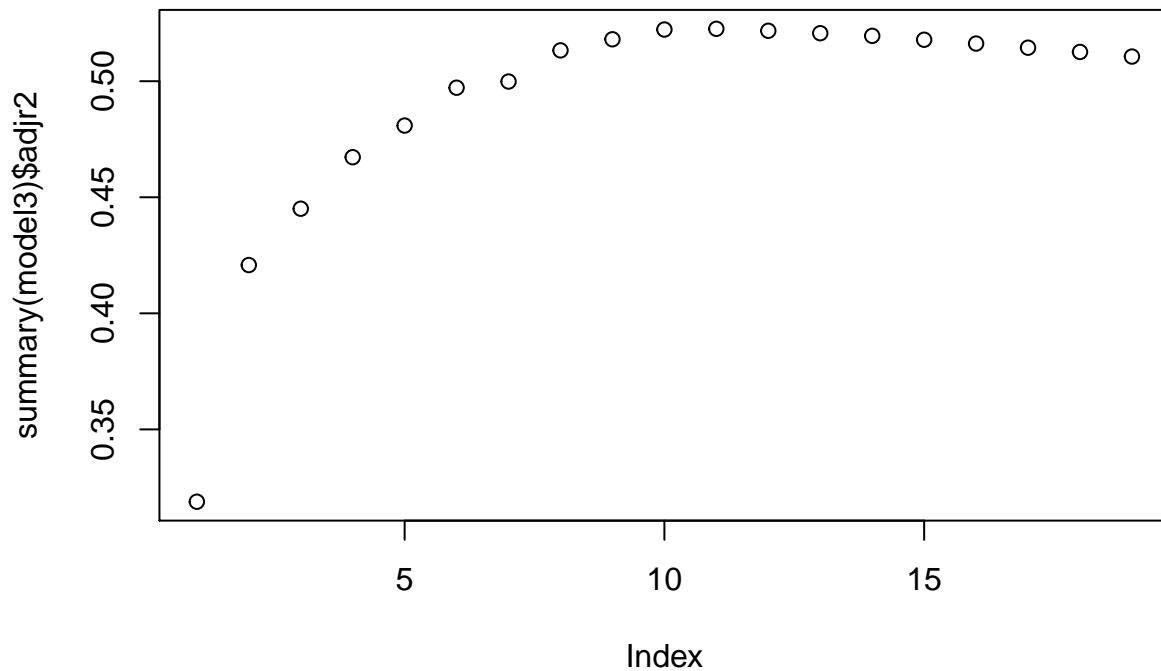
```
# compare between best subsett and forward selection based on r^2 values
summary(model2)$adjr2 - summary(model3)$adjr2
```

```
##  [1] 3.330669e-16 1.110223e-16 0.000000e+00 0.000000e+00 1.110223e-16
##  [6] 0.000000e+00 9.185854e-04 4.314850e-04 1.110223e-16 1.110223e-16
## [11] 1.110223e-16 0.000000e+00 0.000000e+00 2.220446e-16 1.110223e-16
## [16] 0.000000e+00 0.000000e+00 0.000000e+00 0.000000e+00
```

```
plot(summary(model3)$adjr2)
```



```
# backward selection
model4 <- regsubsets(Salary ~ ., data = hitters, nvmax = 19, method = "backward")
which.max(summary(model4)$adjr2)
```

```
## [1] 11
```

```r
coef(model4, 11)
```

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

```r
summary(model4)
```

```
## 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 CRBI
## 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 ) "*"    "*"     "*"      "*"     "*" "*"    "*"     "*"       "*"      "*"      "*"     "*"
##           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 ) "*"    "*"     "*"      "*"     "*"     "*"    "*"
```