

Lab4

Dr. Purna Gamage

```
library(ISLR)
library(leaps)

library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --
## v ggplot2 3.3.5    v purrr  0.3.4
## v tibble  3.1.6    v dplyr  1.0.7
## v tidyr   1.1.4    v stringr 1.4.0
## v readr   2.1.1    v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

library(caret)

## Loading required package: lattice

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
## lift
```

Textbook Example

```
data(Hitters)
summary(Hitters)
```

##	AtBat	Hits	HmRun	Runs
##	Min. : 16.0	Min. : 1	Min. : 0.00	Min. : 0.00
##	1st Qu.:255.2	1st Qu.: 64	1st Qu.: 4.00	1st Qu.: 30.25
##	Median :379.5	Median : 96	Median : 8.00	Median : 48.00
##	Mean :380.9	Mean :101	Mean :10.77	Mean : 50.91
##	3rd Qu.:512.0	3rd Qu.:137	3rd Qu.:16.00	3rd Qu.: 69.00
##	Max. :687.0	Max. :238	Max. :40.00	Max. :130.00
##				
##	RBI	Walks	Years	CAtBat
##	Min. : 0.00	Min. : 0.00	Min. : 1.000	Min. : 19.0
##	1st Qu.: 28.00	1st Qu.: 22.00	1st Qu.: 4.000	1st Qu.: 816.8
##	Median : 44.00	Median : 35.00	Median : 6.000	Median : 1928.0
##	Mean : 48.03	Mean : 38.74	Mean : 7.444	Mean : 2648.7
##	3rd Qu.: 64.75	3rd Qu.: 53.00	3rd Qu.:11.000	3rd Qu.: 3924.2
##	Max. :121.00	Max. :105.00	Max. :24.000	Max. :14053.0

```
##
##      CHits      CHmRun      CRuns      CRBI
## Min.   : 4.0    Min.   : 0.00   Min.   : 1.0    Min.   : 0.00
## 1st Qu.: 209.0  1st Qu.: 14.00   1st Qu.: 100.2  1st Qu.: 88.75
## Median : 508.0  Median : 37.50   Median : 247.0  Median : 220.50
## Mean   : 717.6  Mean   : 69.49   Mean   : 358.8  Mean   : 330.12
## 3rd Qu.:1059.2  3rd Qu.: 90.00   3rd Qu.: 526.2  3rd Qu.: 426.25
## Max.   :4256.0  Max.   :548.00   Max.   :2165.0  Max.   :1659.00
##
##      CWalks      League Division PutOuts      Assists
## Min.   : 0.00   A:175   E:157   Min.   : 0.0    Min.   : 0.0
## 1st Qu.: 67.25  N:147   W:165   1st Qu.: 109.2  1st Qu.: 7.0
## Median : 170.50                Median : 212.0  Median : 39.5
## Mean   : 260.24                Mean   : 288.9  Mean   :106.9
## 3rd Qu.: 339.25                3rd Qu.: 325.0  3rd Qu.:166.0
## Max.   :1566.00                Max.   :1378.0  Max.   :492.0
##
##      Errors      Salary      NewLeague
## Min.   : 0.00   Min.   : 67.5   A:176
## 1st Qu.: 3.00   1st Qu.: 190.0  N:146
## Median : 6.00   Median : 425.0
## Mean   : 8.04   Mean   : 535.9
## 3rd Qu.:11.00   3rd Qu.: 750.0
## Max.   :32.00   Max.   :2460.0
##
##              NA's      :59
```

Need to remove missing values.

```
# remove rows with missing data
Hitters=na.omit(Hitters)
with(Hitters, sum(is.na(Hitters$Salary)))
```

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

Best subset

```
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 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 ) "*" " " "*" "*" " " " " "

```

looks like it only goes up to 8 best variables.

Let's try all 19 variables.

```

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

```

```

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

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

```
reg.summary$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
```

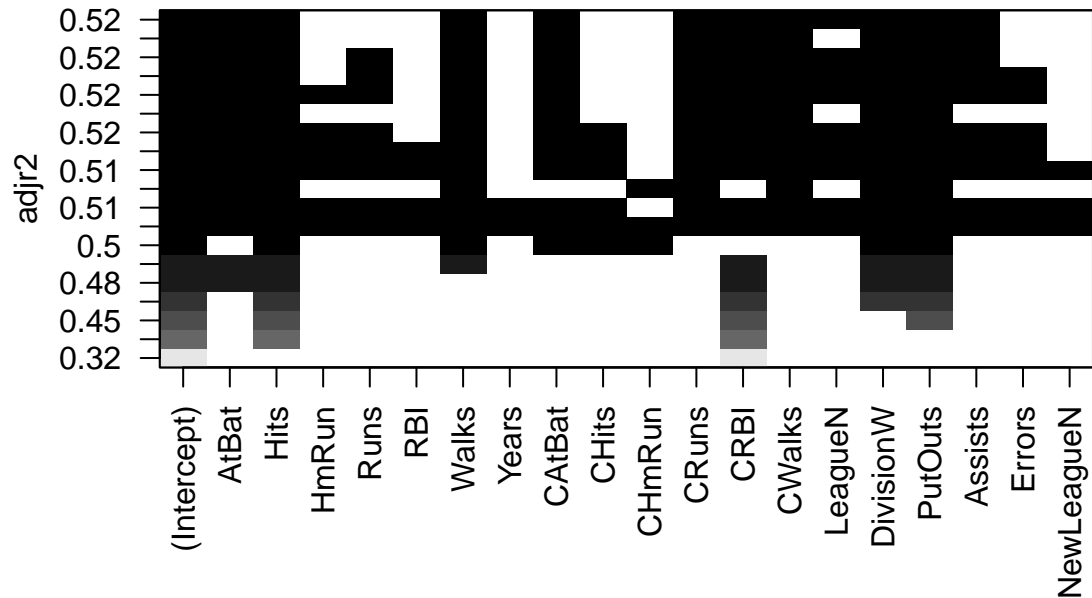
```
summary(regfit.full)$adjr2
```

```
## [1] 0.3188503 0.4208024 0.4450753 0.4672734 0.4808971 0.4972001 0.5007849
## [8] 0.5137083 0.5180572 0.5222606 0.5225706 0.5217245 0.5206736 0.5195431
## [15] 0.5178661 0.5162219 0.5144464 0.5126097 0.5106270
```

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

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

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



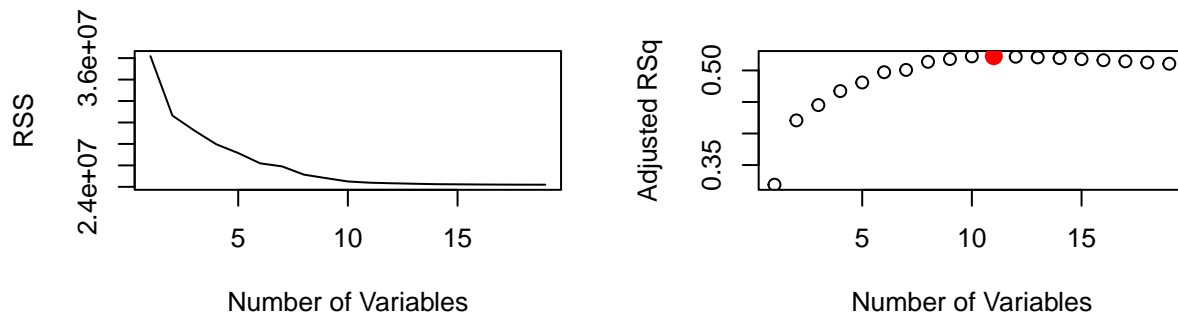
```
par(mfrow=c(2,2))
```

```
plot(reg.summary$rss,xlab="Number of Variables",ylab="RSS",type="l")
```

```
plot(reg.summary$adjr2,xlab="Number of Variables",ylab="Adjusted RSq")
```

```
m=which.max(reg.summary$adjr2)
```

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

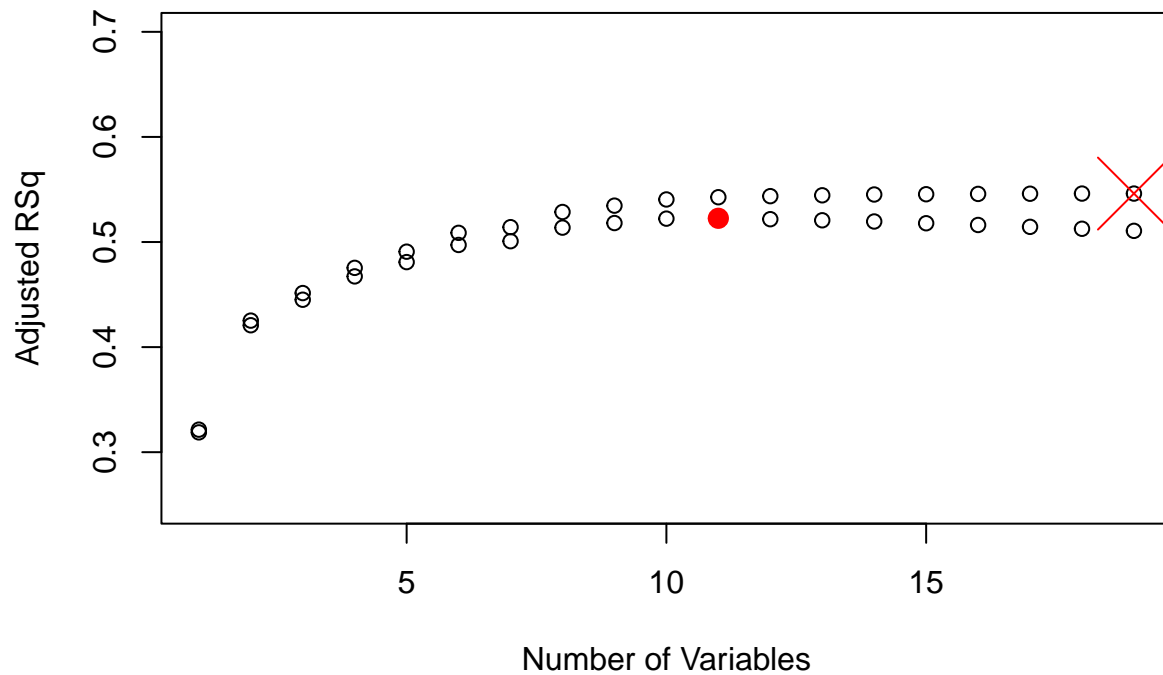


```
#compare Rsq and adjusted srq
```

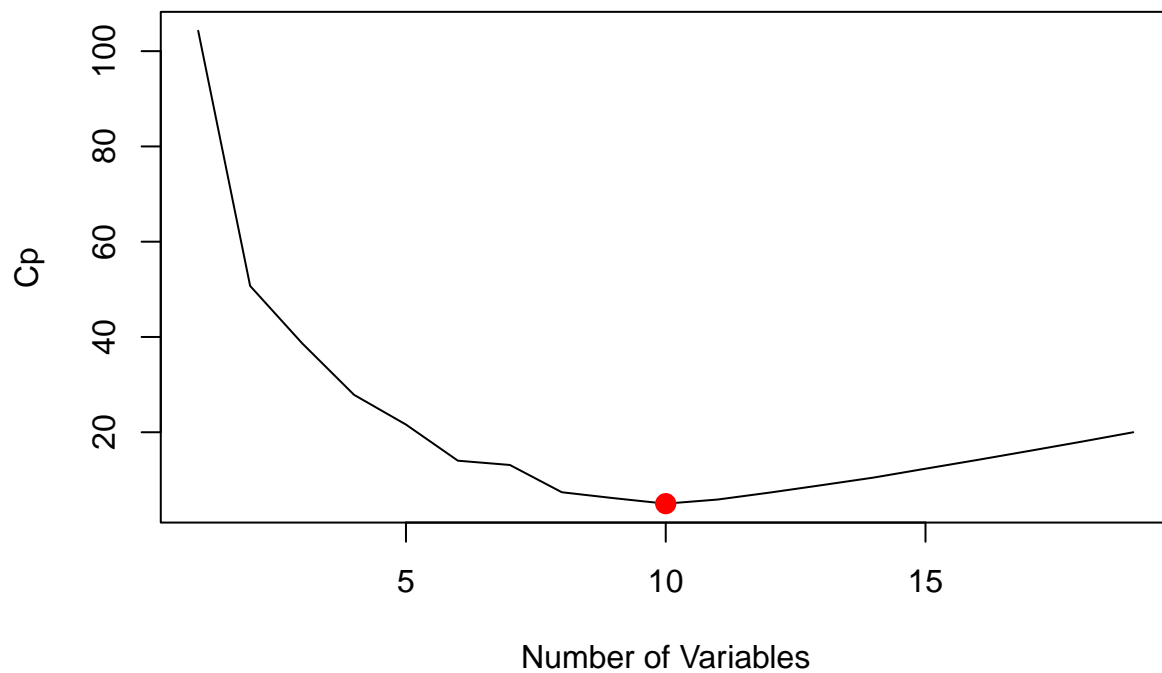
```
plot(reg.summary$adjr2,xlab="Number of Variables",ylab="Adjusted RSq", ylim=c(0.25,0.7))
```

```
m=which.max(reg.summary$adjr2)
```

```
n=which.max(reg.summary$rsq)
points(m,reg.summary$adjr2[m], col="red",cex=2,pch=20)
points(reg.summary$rsq)
points(n,reg.summary$rsq[n], col="red",cex=5,pch=4)
```

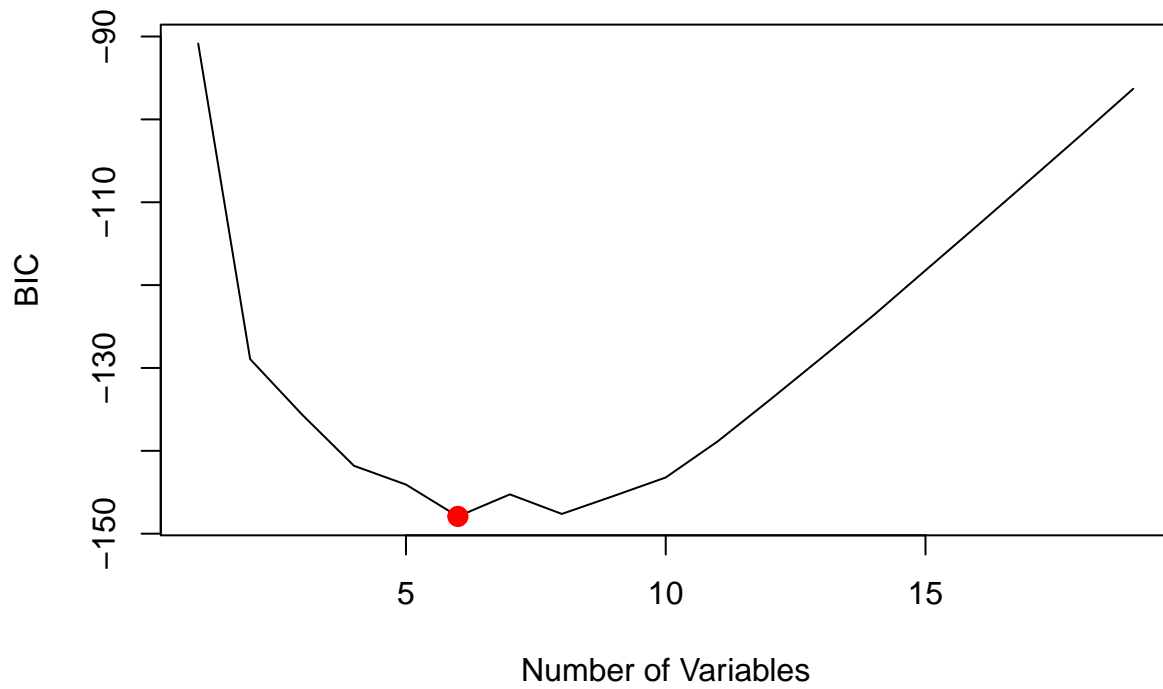


```
plot(reg.summary$cp,xlab="Number of Variables",ylab="Cp",type='l')
m=which.min(reg.summary$cp)
points(m,reg.summary$cp[m], col="red",cex=2,pch=20)
```

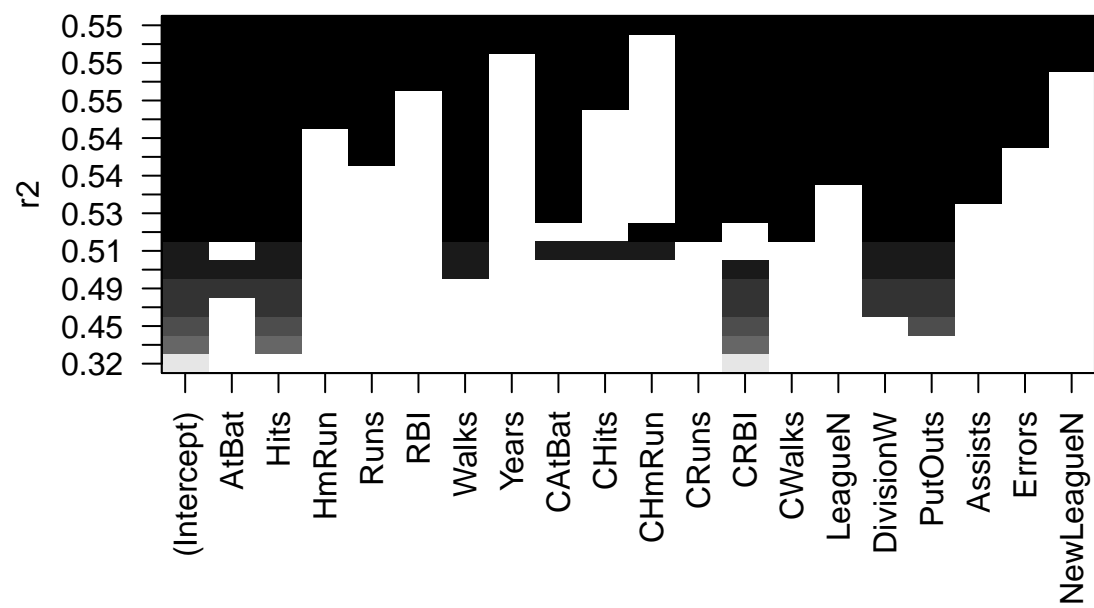


```
m=which.min(reg.summary$bic)
plot(reg.summary$bic,xlab="Number of Variables",ylab="BIC",type='l')
```

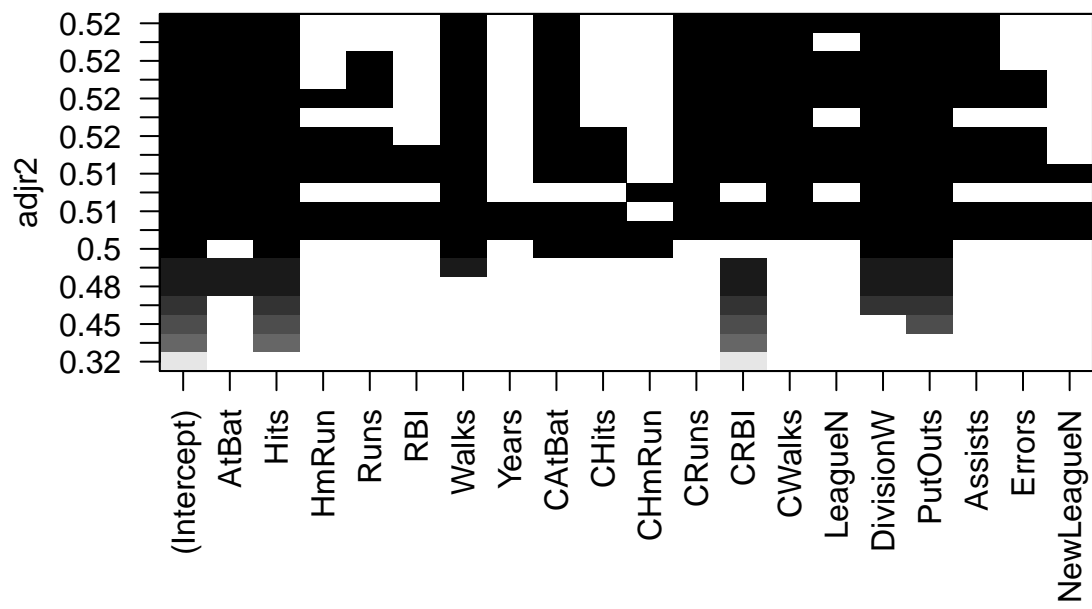
```
points(m,reg.summary$bic[6],col="red",cex=2,pch=20)
```



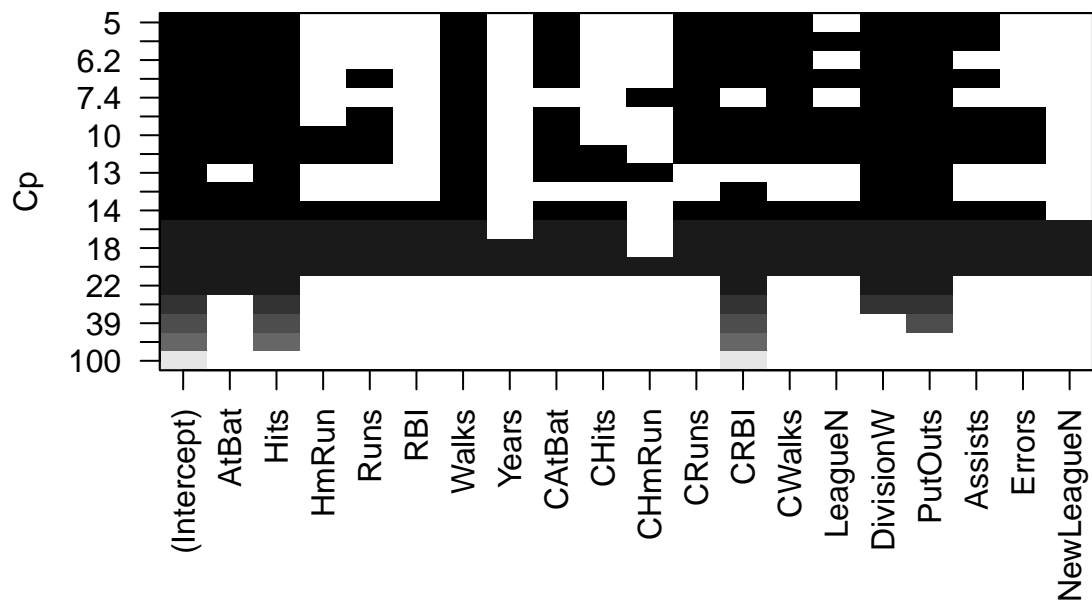
```
plot(regfit.full,scale="r2")
```



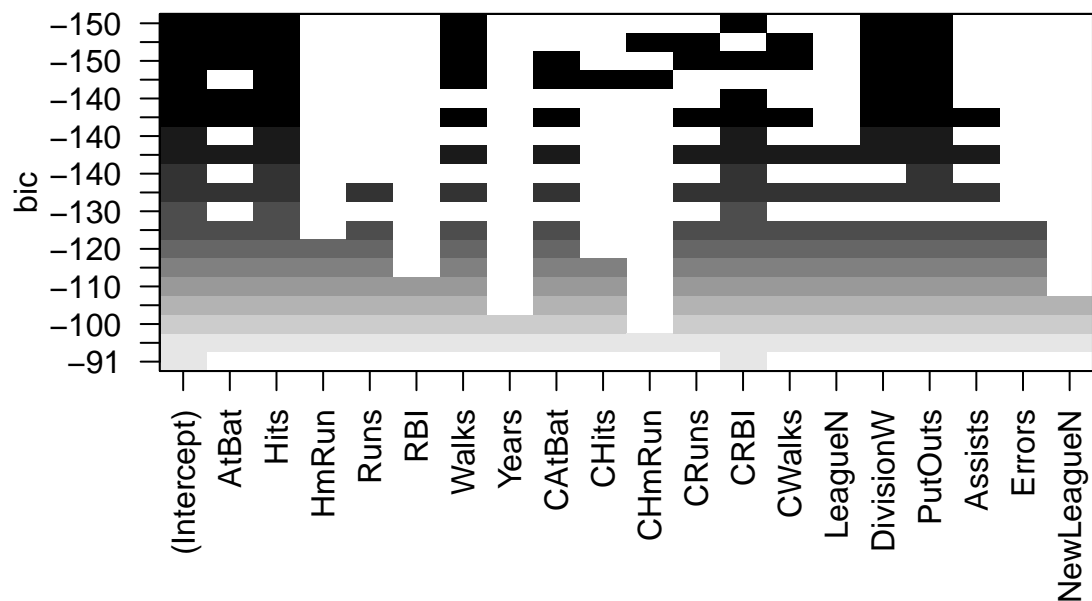
```
plot(regfit.full,scale="adjr2")
```



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



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

```
coef(regfit.full,6)
```

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

Forward and Backward Stepwise Selection

```
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 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 ) "*" "*" " "*" "*" "*" "*" "*" "*" " "

```

variables are nested.

```

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	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)	"*	" "	"*	"*	"*	" "	" "						

```
## 15 ( 1 ) "*" "*" "*" "*" "*" " "
## 16 ( 1 ) "*" "*" "*" "*" "*" "*" " "
## 17 ( 1 ) "*" "*" "*" "*" "*" "*" "*"
## 18 ( 1 ) "*" "*" "*" "*" "*" "*" "*"
## 19 ( 1 ) "*" "*" "*" "*" "*" "*" "*"

```

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

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

```

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

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

```

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

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

```

Hybrid method

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

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

```

Choosing Among Models

In order to use the validation set approach, we begin by splitting the observations into a training set and a test set.

We do this by creating a random vector, `train`, of elements equal to `TRUE` if the corresponding observation is in the training set, and `FALSE` otherwise.

The vector `test` has a `TRUE` if the observation is in the test set, and a `FALSE` otherwise.

Note the `!` in the command to create `test` causes `TRUE`s to be switched to `FALSE`s and vice versa.

We also set a random seed so that the user will obtain the same training set/test set split.

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

Now, we apply `regsubsets()` to the training set in order to perform best subset selection.

Notice that we subset the Hitters data frame directly in the call in order to access only the training subset of the data, using the expression `Hitters[train,]`.

```
regfit.best=regsubsets(Salary~.,data=Hitters[train,],nvmax=19)
```

We now compute the validation set error for the best model of each model size.

We first make a model matrix from the test data.

The `model.matrix()` function is used in many regression packages for building an “X” matrix from data.

```
test.mat=model.matrix(Salary~.,data=Hitters[test,])
```

Now we run a loop, and for each size `i`, we extract the coefficients from `regfit.best` for the best model of that size, multiply them into the appropriate columns of the test model matrix to form the predictions, and compute the test MSE.

```
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)
}
```

We find that the best model is the one that contains ten variables.

```
val.errors
```

```
## [1] 164377.3 144405.5 152175.7 145198.4 137902.1 139175.7 126849.0 136191.4
## [9] 132889.6 135434.9 136963.3 140694.9 140690.9 141951.2 141508.2 142164.4
## [17] 141767.4 142339.6 142238.2
```

```
m=which.min(val.errors) #m=10
coef(regfit.best,m)
```

```
## (Intercept)      AtBat      Hits      Walks      CRuns      CWalks
##  67.1085369  -2.1462987   7.0149547   8.0716640   1.2425113  -0.8337844
##   DivisionW      PutOuts
## -118.4364998   0.2526925
```

This was a little tedious, partly because there is no `predict()` method for `regsubsets()`.

Since we will be using this function again, we can capture our steps above and write our own `predict` method.

```
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
}
```

Our function pretty much mimics what we did above.

Finally, we perform best subset selection on the full data set, and select the best ten-variable model.

```
regfit.best=regsubsets(Salary~.,data=Hitters,nvmax=19)
coef(regfit.best,10)
```

```
## (Intercept)      AtBat      Hits      Walks      CAtBat      CRuns
## 162.5354420    -2.1686501    6.9180175    5.7732246   -0.1300798    1.4082490
##          CRBI      CWalks    DivisionW    PutOuts      Assists
##      0.7743122   -0.8308264  -112.3800575    0.2973726    0.2831680
```

In fact, we see that the best ten-variable model on the full data set has a different set of variables than the best ten-variable model on the training set.

We now try to choose among the models of different sizes using cross- validation.

we must perform best subset selection within each of the k training sets.

Despite this, we see that with its clever subsetting syntax, R makes this job quite easy.

First, we create a vector that allocates each observation to one of $k = 10$ folds, and we create a matrix in which we will store the results.

```
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)))
```

Now we write a for loop that performs cross-validation.

In the jth fold, the elements of folds that equal j are in the test set, and the remainder are in the training set.

We make our predictions for each model size (using our new predict() method), compute the test errors on the appropriate subset, and store them in the appropriate slot in the matrix cv.errors.

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

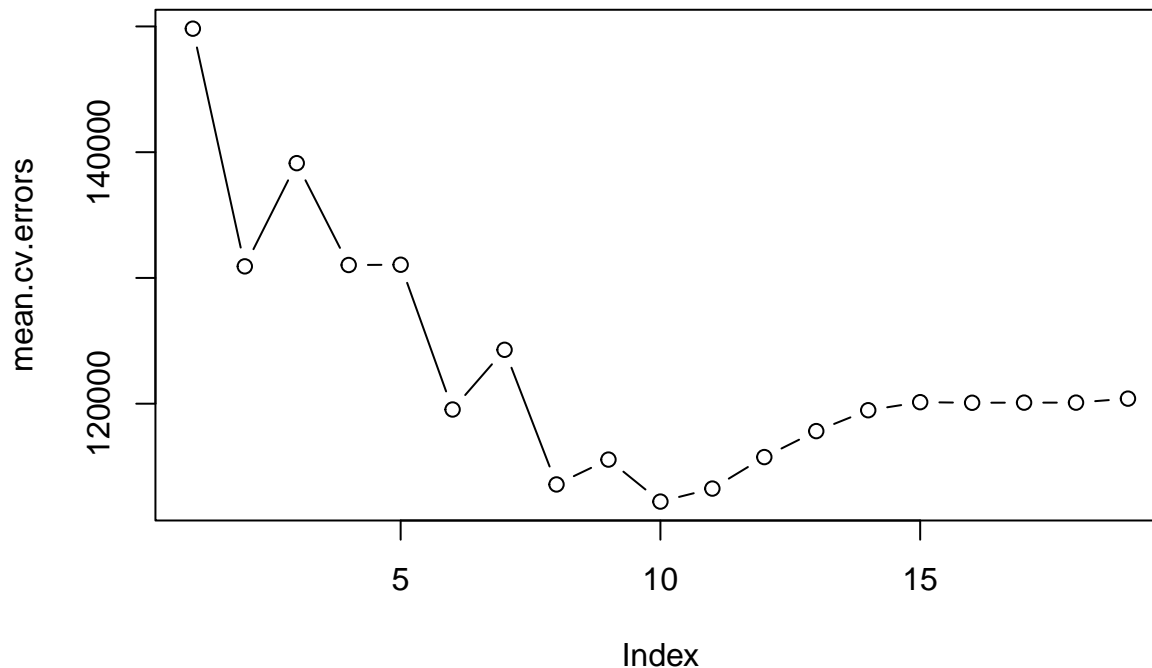
This has given us a 10×19 matrix, of which the (i, j)th element corresponds to the test MSE for the ith cross-validation fold for the best j-variable model.

We use the apply() function to average over the columns of this matrix in order to obtain a vector for which the jth element is the cross- validation error for the j-variable model.

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

```
##      1      2      3      4      5      6      7      8
## 149821.1 130922.0 139127.0 131028.8 131050.2 119538.6 124286.1 113580.0
##      9     10     11     12     13     14     15     16
## 115556.5 112216.7 113251.2 115755.9 117820.8 119481.2 120121.6 120074.3
##     17     18     19
## 120084.8 120085.8 120403.5
```

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



We see that cross-validation selects an 11-variable model.

We now perform best subset selection on the full data set in order to obtain the 11-variable model.

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

Example

Predicting fertility score on the basis of socio-economic indicators.

```
# Load the data
data("swiss")
# Inspect the data
sample_n(swiss, 3)
```

```
##      Fertility Agriculture Examination Education Catholic Infant.Mortality
## Nyone      56.6      50.9      22      12      15.14      16.7
## Moutier     85.8      36.5      12      7      33.77      20.3
## Echallens   68.3      72.6      18      2      24.20      21.2
```

(a). Computing best subsets regression.

In our example, we have only 5 predictor variables in the data. So, we'll use `nvmax = 5`.

```
models <- regsubsets(Fertility~., data = swiss, nvmax = 5)
summary(models)
```

```
## Subset selection object
## Call: regsubsets.formula(Fertility ~ ., data = swiss, nvmax = 5)
## 5 Variables (and intercept)
```



```
##              Forced in Forced out
## Agriculture      FALSE      FALSE
## Examination      FALSE      FALSE
## Education        FALSE      FALSE
## Catholic         FALSE      FALSE
## Infant.Mortality FALSE      FALSE
## 1 subsets of each size up to 5
## Selection Algorithm: exhaustive
##      Agriculture Examination Education Catholic Infant.Mortality
## 1  ( 1 ) " "      " "      "*"      " "      " "
## 2  ( 1 ) " "      " "      "*"      "*"      " "
## 3  ( 1 ) " "      " "      "*"      "*"      "*"
## 4  ( 1 ) "*"      " "      "*"      "*"      "*"
## 5  ( 1 ) "*"      "*"      "*"      "*"      "*"

```

It can be seen that the best 2-variables model contains only Education and Catholic variables ($Fertility \sim Education + Catholic$).

The best three-variable model is

($Fertility \sim Education + Catholic + Infant.mortality$), and so forth.

A natural question is: which of these best models should we finally choose for our predictive analytics?

(c) Choosing the optimal model.

```
res.sum <- summary(models)

data.frame(
  Adj.R2 = which.max(res.sum$adjr2),
  CP = which.min(res.sum$cp),
  BIC = which.min(res.sum$bic)
)
```

```
##   Adj.R2 CP BIC
## 1     5  4  4

```

Here, adjusted R2 tells us that the best model is the one with all the 5 predictor variables. However, using the BIC and Cp criteria, we should go for the model with 4 variables.

Note also that the adjusted R2, BIC and Cp are calculated on the training data that have been used to fit the model. This means that, the model selection, using these metrics, is possibly subject to overfitting and may not perform as well when applied to new data.

A more rigorous approach is to select a models based on the prediction error computed on a new test data using k-fold cross-validation

(d). K-fold cross-validation

(i) `get_model_formula()`, allowing to access easily the formula of the models returned by the function `regsubsets()`.

```
# id: model id
# object: regsubsets object
# data: data used to fit regsubsets
# outcome: outcome variable
get_model_formula <- function(id, object, outcome){
  # get models data
  models <- summary(object)$which[id,-1]
  # Get outcome variable

```

```

#form <- as.formula(object$call[[2]])
#outcome <- all.vars(form)[1]
# Get model predictors
predictors <- names(which(models == TRUE))
predictors <- paste(predictors, collapse = "+")
# Build model formula
as.formula(paste0(outcome, "~", predictors))
}

```

For example to have the best 3-variable model formula;

```
get_model_formula(3, models, "Fertility")
```

```
## Fertility ~ Education + Catholic + Infant.Mortality
## <environment: 0x7fd8080cbbd8>
```

(ii) `get_cv_error()`, to get the cross-validation (CV) error for a given model:

```

get_cv_error <- function(model.formula, data){
  set.seed(1)
  train.control <- trainControl(method = "cv", number = 5)
  cv <- train(model.formula, data = data, method = "lm",
             trControl = train.control)
  cv$results$RMSE
}

```

Finally, use the above defined helper functions to compute the prediction error of the different best models returned by the `regsubsets()` function:

```

# Compute cross-validation error
model.ids <- 1:5
cv.errors <- map(model.ids, get_model_formula, models, "Fertility") %>%
  map(get_cv_error, data = swiss) %>%
  unlist()
cv.errors

```

```
## [1] 9.464156 8.517433 7.855267 7.601072 7.736328
```

```

# Select the model that minimize the CV error
which.min(cv.errors)

```

```
## [1] 4
```

It can be seen that the model with 4 variables is the best model. It has the lower prediction error.

The regression coefficients of this model can be extracted as follow:

```
coef(models, 4)
```

```

##      (Intercept)      Agriculture      Education      Catholic
##      62.1013116      -0.1546175      -0.9802638      0.1246664
## Infant.Mortality
##      1.0784422

```

Read more:

<http://www.sthda.com/english/articles/37-model-selection-essentials-in-r/155-best-subsets-regression-essentials-in-r/>

<http://www.sthda.com/english/articles/37-model-selection-essentials-in-r/154-stepwise-regression-essentials-in-r/>

Lab Assignment

Explore the a different dataset from the ISLR2 package (ISLR2::) and use Best Subset selection with the `regsubsets()` function. From the outputs of `regsubsets()`, we can obtain metrics like BIC for each sub-model with K predictors.

1. Make a few plots showing BIC, and Mallows's Cp, and 1-AdjustedRsquared.

Are the different metrics in rough agreement?

Try to answer these questions using different plots

- a. compare RSS and R^2
 - b. Compare R^2 with adjusted R^2
 - c. Compare adjusted R^2 , BIC, cp with Cross validation results
 - d. Also get the feature plots for Adj. R^2 , BIC, cp; What is the number of predictors to include in the model for each different metric and what are these predictors?
2. These metrics are evaluted with the full training set. Another way we can approach the problem of model selection is to use cross validation. What are the pros and cons of each approach?