```
Cross-validation and Model Selection
Lab: Model Selection
Subset Selection Methods
Best Subset Selection
Here we apply the best subset selection approach to the Hitters data. We wish to predict a baseball player's Salary on the basis of various
statistics associated with performance in the previous year.
First of all, we note that the Salary variable is missing for some of the players. The is.na() function can be used to identify the missing
observations. It returns a vector of the same length as the input vector, with a TRUE for any elements that are missing, and a FALSE for non-
missing elements. The sum() function can then be used to count all of the missing elements.
 library(ISLR2)
 ## Warning: package 'ISLR2' was built under R version 4.0.5
 names(Hitters)
                                                            "RBI"
 ## [1] "AtBat"
                      "Hits"
                                   "HmRun"
                                               "Runs"
                                                                         "Walks"
                      "CAtBat"
 ## [7] "Years"
                                   "CHits"
                                               "CHmRun"
                                                           "CRuns"
                                                                        "CRBI"
 ## [13] "CWalks"
                     "League"
                                  "Division" "PutOuts" "Assists"
                                                                        "Errors"
                      "NewLeague"
 ## [19] "Salary"
 dim(Hitters)
 ## [1] 322 20
 sum(is.na(Hitters$Salary))
 ## [1] 59
Hence we see that Salary is missing for 59 players. The na.omit() function removes all of the rows that have missing values in any variable.
 Hitters <- na.omit(Hitters)</pre>
 dim(Hitters)
 ## [1] 263 20
 sum(is.na(Hitters))
 ## [1] 0
The regsubsets() function (part of the leaps library) performs best subset selection by identifying the best model that contains a given number
of predictors, where best is quantified using RSS. The syntax is the same as for lm(). The summary() command outputs the best set of
variables for each model size.
 library(leaps)
 regfit.full <- regsubsets(Salary ~ ., Hitters)</pre>
 summary(regfit.full)
 ## Subset selection object
 ## Call: regsubsets.formula(Salary ~ ., Hitters)
 ## 19 Variables (and intercept)
                Forced in Forced out
                    FALSE
 ## AtBat
                    FALSE
                               FALSE
 ## Hits
 ## HmRun
                    FALSE
 ## Runs
                    FALSE
                               FALSE
 ## RBI
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## Walks
 ## Years
                    FALSE
                               FALSE
 ## CAtBat
                    FALSE
                               FALSE
                    FALSE
 ## CHits
                               FALSE
                    FALSE
                               FALSE
 ## CHmRun
                    FALSE
                               FALSE
 ## CRuns
 ## CRBI
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## CWalks
 ## LeagueN
                    FALSE
                               FALSE
                               FALSE
                    FALSE
 ## DivisionW
                    FALSE
                               FALSE
 ## PutOuts
 ## Assists
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## Errors
 ## NewLeagueN
                   FALSE
 ## 1 subsets of each size up to 8
 ## Selection Algorithm: exhaustive
              AtBat Hits HmRun Runs RBI Walks Years CAtBat CHits CHmRun CRuns CRBI
 CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
 ## 1 ( 1 ) " " " "
 ## 2 ( 1 ) " " " " "
 ## 3 (1) "" ""
 ## 4 ( 1 ) " " " "
                                    11 * 11
 ## 5 ( 1 ) " " " "
                                    11 * 11
 ## 6 ( 1 ) " " " "
 ## 7 (1)""""*"
                                   II * II
 ## 8 ( 1 ) "*" " " "*"
An asterisk indicates that a given variable is included in the corresponding model. For instance, this output indicates that the best two-variable
model contains only Hits and CRBI. By default, regsubsets() only reports results up to the best eight-variable model. But the nvmax option
can be used in order to return as many variables as are desired. Here we fit up to a 19-variable model.
 regfit.full <- regsubsets(Salary ~ ., data = Hitters,</pre>
     nvmax = 19)
 reg.summary <- summary(regfit.full)</pre>
The summary() function also returns \mathbb{R}^2, RSS, adjusted \mathbb{R}^2, \mathbb{C}_p, and BIC. We can examine these to try to select the best overall model.
 names(reg.summary)
                                                       "bic"
 ## [1] "which" "rsq"
                           "rss"
                                    "adjr2" "cp"
                                                                "outmat" "obj"
For instance, we see that the R^2 statistic increases from 32\,\%, when only one variable is included in the model, to almost 55\,\%, when all
variables are included. As expected, the \mathbb{R}^2 statistic increases monotonically as more variables are included.
 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
Plotting RSS, adjusted R^2, C_p, and BIC for all of the models at once will help us decide which model to select. Note the type = "1" option tells
R to connect the plotted points with lines.
 par(mfrow = c(1, 2))
 plot(reg.summary$rss, xlab = "Number of Variables",
     ylab = "RSS", type = "1")
 plot(reg.summary$adjr2, xlab = "Number of Variables",
     ylab = "Adjusted RSq", type = "1")
     3.6e+07
                                                   0.50
      3.2e+07
                                                   0.45
                                             Adjusted RSq
                                                   0.40
      2.8e+07
                                                   0.35
      2.4e+07
                5
                        10
                               15
                                                                     10
                                                                            15
               Number of Variables
                                                            Number of Variables
The points() command works like the plot() command, except that it puts points on a plot that has already been created, instead of creating
a new plot. The which.max() function can be used to identify the location of the maximum point of a vector. We will now plot a red dot to indicate
the model with the largest adjusted \mathbb{R}^2 statistic.
 which.max(reg.summary$adjr2)
 ## [1] 11
 plot(reg.summary$adjr2, xlab = "Number of Variables",
     ylab = "Adjusted RSq", type = "1")
 points(11, reg.summary$adjr2[11], col = "red", cex = 2,
     pch = 20)
      50
     0.45
Adjusted RSq
     0.40
      0.35
                                              10
                                                                 15
                                      Number of Variables
In a similar fashion we can plot the C_p and BIC statistics, and indicate the models with the smallest statistic using which.min().
 plot(reg.summary$cp, xlab = "Number of Variables",
     ylab = "Cp", type = "1")
 which.min(reg.summary$cp)
 ## [1] 10
 points(10, reg.summaryscp[10], col = "red", cex = 2,
     pch = 20)
      100
      80
      9
 _{0}^{C}
      40
      20
                                              10
                                                                 15
                                      Number of Variables
 which.min(reg.summary$bic)
 ## [1] 6
 plot(reg.summary$bic, xlab = "Number of Variables",
     ylab = "BIC", type = "1")
 points(6, reg.summary$bic[6], col = "red", cex = 2,
     pch = 20)
      90
      -100
      -110
      -120
      -130
      -140
      -150
                                              10
                                                                 15
                                      Number of Variables
The regsubsets() function has a built-in plot() command which can be used to display the selected variables for the best model with a given
number of predictors, ranked according to the BIC, C_p, adjusted R^2, or AIC. To find out more about this function, type <code>?plot.regsubsets</code> .
 plot(regfit.full, scale = "r2")
      0.55
      0.55
      0.55
      0.54
  인 0.54
      0.53
      0.51
      0.49
      0.45
      0.32
                        HmRun
Runs
RBI
                                  Walks
Years
CAtBat
                                                CHmRun
CRuns
                                                              LeagueN
DivisionW
                                             CHits
                                                           CWalks
                                                                     PutOuts
                                                       CRBI
                                                                        Assists
                                                                               NewLeagueN
              (Intercept)
 plot(regfit.full, scale = "adjr2")
      0.52
      0.52
      0.52
      0.52
   adjr2
0.51
       0.5
      0.48
      0.32
                                                              LeagueN
                 AtBat
                     Hits
                        HmRun
                            Runs
                               RBI
                                  Walks
                                      Years
                                         CAtBat
                                             CHits
                                                CHmRun
                                                   CRuns
                                                       CRBI
                                                          CWalks
                                                                     PutOuts
                                                                         Assists
                                                                            Errors
                                                                               NewLeagueN
              (Intercept)
                                                                 DivisionW
 plot(regfit.full, scale = "Cp")
       6.2
       7.4
        10
        13
   Cp
        18
        22
        39
       100
                        HmRun
                 AtBat
                                                                     PutOuts
                     Hits
                           Runs
                                  Walks
                                         CAtBat
                                            CHits
                                                CHmRun
                                                                         Assists
                               RBI
                                      Years
                                                    CRuns
                                                       CRBI
                                                          CWalks
                                                                                NewLeagueN
              (Intercept)
                                                              LeagueN
                                                                 DivisionW
 plot(regfit.full, scale = "bic")
      -150
      -150
      -140
      -140
      -130
      -120
      -110
      -100
       -91
                                                CHmRun
CRuns
CRBI
                                                             LeagueN
DivisionW
PutOuts
                                  Walks
Years
CAtBat
CHits
                        HmRun
Runs
                                                          CWalks
                               RBI
                                                                                NewLeagueN
              (Intercept)
The top row of each plot contains a black square for each variable selected according to the optimal model associated with that statistic. For
instance, we see that several models share a BIC close to -150. However, the model with the lowest BIC is the six-variable model that contains
only AtBat, Hits, Walks, CRBI, DivisionW, and PutOuts. We can use the coef() function to see the coefficient estimates associated
with this model.
 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
We can also use the regsubsets() function to perform forward stepwise or backward stepwise selection, using the argument
method = "forward" or method = "backward".
 regfit.fwd <- regsubsets(Salary ~ ., data = Hitters,</pre>
     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
                    FALSE
                               FALSE
 ## Hits
 ## HmRun
                    FALSE
                               FALSE
 ## Runs
                    FALSE
                               FALSE
 ## RBI
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## Walks
 ## Years
                    FALSE
                               FALSE
 ## CAtBat
                    FALSE
                               FALSE
 ## CHits
                    FALSE
                               FALSE
 ## CHmRun
                    FALSE
                               FALSE
 ## CRuns
                    FALSE
                               FALSE
 ## CRBI
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## CWalks
                    FALSE
                               FALSE
 ## LeagueN
 ## DivisionW
                    FALSE
                               FALSE
 ## PutOuts
                    FALSE
                               FALSE
                    FALSE
 ## Assists
                               FALSE
 ## Errors
                    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) "" "" "" "" "" ""
 ## 3 (1) "" "*" "" "" "" "" "" ""
 ## 4 ( 1 )
 ## 5 (1) "*" "*" " " " " " " " " " " "
 ## 6 (1) "*" "*" " " " " " " " " " " " "
              11 * 11
 11 * 11
        ## 14 ( 1 ) "*" "*" "*" "*" """ ""
 ## 15 ( 1 ) "*" "*" "*" "*" "*" "*"
                                                              11 * 11 11
                                                                           11 * 11 * 11 * 11
 ## 16 ( 1 ) "*" "*" "*"
 ## 17 ( 1 ) "*" "*" "*"
 ## 18 ( 1 ) "*" "*" "*" "*" "*" "*" "*"
                                                             11 * 11 11
                                                                           11 * 11
 ## 19 ( 1 ) "*" "*" "*" "*" "*" "*"
               CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
 ## 1 ( 1 ) " " " "
 ## 2 ( 1 ) " "
                                        11 * 11
 ## 3 (1)
 ## 4 ( 1 ) " " " "
                              11 * 11
                                        11 * 11
                                                         11 11
 ## 5 (1) """ ""
                              II * II
                                        11 * 11
                                                         ш
                              II * II
                                        II * II
 ## 6 (1)
 ## 7 ( 1 ) "*" " "
                              11 * 11
                                        11 * 11
 ## 8 ( 1 ) "*"
                              II * II
                                         II * II
 ## 9 ( 1 ) "*"
                                        II * II
 ## 10 ( 1 ) "*" " "
                                        II * II
 ## 11 ( 1 ) "*"
                              II * II
                                         11 * 11
                                                 11 * 11
                                                         11 11
 ## 12 ( 1 ) "*"
                                         II * II
 ## 13 ( 1 ) "*"
                                         II * II
 ## 14 ( 1 ) "*"
                              11 * 11
                                         II * II
                                                 11 * 11
                                                         11 * 11
                                                         11 * 11
 ## 15 ( 1 ) "*"
 ## 16 ( 1 ) "*"
                                         11 * 11
                                         II * II
 ## 17 ( 1 ) "*"
                              11 * 11
                                                 11 * 11
                                                         11 * 11
 ## 18 ( 1 ) "*"
                                                         11 * 11
 ## 19 ( 1 ) "*" "*"
                                         11 * 11
 regfit.bwd <- regsubsets(Salary ~ ., data = Hitters,</pre>
     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
                               FALSE
 ## Hits
                    FALSE
                               FALSE
 ## HmRun
 ## Runs
                    FALSE
                    FALSE
 ## RBI
                               FALSE
                    FALSE
                               FALSE
 ## Walks
 ## Years
                    FALSE
                               FALSE
 ## CAtBat
                    FALSE
                                FALSE
 ## CHits
                    FALSE
                                FALSE
 ## CHmRun
                    FALSE
                                FALSE
 ## CRuns
                    FALSE
                               FALSE
 ## CRBI
                    FALSE
                               FALSE
 ## CWalks
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## LeagueN
 ## DivisionW
                    FALSE
                               FALSE
                    FALSE
                               FALSE
 ## PutOuts
                    FALSE
                               FALSE
 ## Assists
 ## 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) ""
 11 * 11
 ## 5 ( 1 )
 ## 6 ( 1 )
               11*11 11*11 11 11 11 11 11 11 11 11 11
                                                                           11 * 11
 ## 7 (1)
               11 * 11 * 11
 ## 9 ( 1 ) "*" "*"
 11 11 11 11
                                                                           11 * 11
 ## 11 ( 1 ) "*"
        (1)
        11 * 11
        ( 1 ) "*" "*" "*"
                                11 * 11
        (1)"*" "*" "*"
 ## 15
 ## 16 ( 1 ) "*" "*" "*" "*" "*" "*"
                                                              11 * 11 | 11
                                                                           11 * 11
 ## 17 ( 1 ) "*" "*" "*"
                                 || * || || * || || * || || || || || || * ||
 ## 18 ( 1 ) "*" "*" "*"
 ## 19 ( 1 ) "*" "*" "*" "*" "*"
               CWalks LeagueN DivisionW PutOuts Assists Errors NewLeagueN
 ## 1 ( 1 )
 ## 2 ( 1 )
                                         11 * 11
 ## 3 ( 1 )
                                        11 * 11
 ## 4 ( 1 )
 ## 5 ( 1 ) " "
                                                 11 11
              11 11
                              11 * 11
                                        11 * 11
                                                         11 11
 ## 6 (1)
                                         11 * 11
              11 * 11
 ## 8 ( 1 ) "*"
                              11 * 11
                                         11 * 11
 ## 9 ( 1 ) "*"
                              11 * 11
                                         11 * 11
                                                         11 11
 ## 10 ( 1 ) "*"
 ## 11 ( 1 ) "*"
                                         11 * 11
                                        11 * 11
 ## 12 ( 1 ) "*"
                              II * II
                                                         11 11
                                                 11 * 11
        ( 1 ) "*"
 ## 13
 ## 14 ( 1 )
                              11 * 11
                                         11 * 11
                                                         11 * 11
 ## 15 ( 1 ) "*"
                              II * II
                                        11 * 11
                                                 II * II
                                                         II * II
                     11 * 11
 ## 16 ( 1 ) "*"
 ## 17 ( 1 ) "*"
                                         11 * 11
                              11 * 11
                                        11 * 11
 ## 18 ( 1 ) "*" "*"
                                                 11 * 11
                                                         II * II
                                                                 11 * 11
 ## 19 ( 1 ) "*" "*"
For instance, we see that using forward stepwise selection, the best one-variable model contains only CRBI, and the best two-variable model
additionally includes Hits. For this data, the best one-variable through six-variable models are each identical for best subset and forward
selection. However, the best seven-variable models identified by forward stepwise selection, backward stepwise selection, and best subset
selection are different.
 coef(regfit.full, 7)
     (Intercept)
                          Hits
 ## 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)
                                       Hits
                                                    Walks
                                                                               CWalks
                         AtBat
                                                                 CRuns
                                  6.7574914 6.0558691 1.1293095 -0.7163346
 ## 105.6487488 -1.9762838
                    PutOuts
 ## DivisionW
 ## -116.1692169 0.3028847
Choosing Among Models Using the Validation-Set Approach and Cross-Validation
We just saw that it is possible to choose among a set of models of different sizes using C_p, BIC, and adjusted R^2. We will now consider how to do
this using the validation set and cross-validation approaches.
In order for these approaches to yield accurate estimates of the test error, we must use only the training observations to perform all aspects of
model-fitting—including variable selection. Therefore, the determination of which model of a given size is best must be made using only the training
observations. This point is subtle but important. If the full data set is used to perform the best subset selection step, the validation set errors and
cross-validation errors that we obtain will not be accurate estimates of the test error.
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),</pre>
     replace = TRUE)
 test <- (!train)</pre>
Now, we apply regsubsets() to the training set in order to perform best subset selection.
 regfit.best <- regsubsets(Salary ~ .,</pre>
     data = Hitters[train, ], nvmax = 19)
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, ]. 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.
 test.mat <- model.matrix(Salary ~ ., data = Hitters[test, ])</pre>
The model.matrix() function is used in many regression packages for building an X matrix from data. 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)</pre>
 for (i in 1:19) {
  coefi <- coef(regfit.best, id = i)</pre>
  pred <- test.mat[, names(coefi)] %*% coefi</pre>
  val.errors[i] <- mean((Hitters$Salary[test] - pred)^2)</pre>
We find that the best model is the one that contains seven variables.
 val.errors
 ## [1] 164377.3 144405.5 152175.7 145198.4 137902.1 139175.7 126849.0 136191.4
 ## [9] 132889.6 135434.9 136963.3 140694.9 140690.9 141951.2 141508.2 142164.4
 ## [17] 141767.4 142339.6 142238.2
 which.min(val.errors)
 ## [1] 7
 coef(regfit.best, 7)
 ## (Intercept)
                         AtBat
                                       Hits
                                                    Walks
                                                                 CRuns
                                                                              CWalks
 ## 67.1085369 -2.1462987 7.0149547 8.0716640 1.2425113 -0.8337844
 ## DivisionW
                     PutOuts
 ## -118.4364998
                     0.2526925
This was a little tedious, partly because there is no predict() method for regsubsets(). Since we will be using this function again, we can
capture our steps above and write our own predict method.
  predict.regsubsets <- function(object, newdata, id, ...) {</pre>
   form <- as.formula(object$call[[2]])</pre>
   mat <- model.matrix(form, newdata)</pre>
   coefi <- coef(object, id = id)</pre>
   xvars <- names(coefi)</pre>
   mat[, xvars] %*% coefi
Our function pretty much mimics what we did above. The only complex part is how we extracted the formula used in the call to regsubsets().
We demonstrate how we use this function below, when we do cross-validation.
Finally, we perform best subset selection on the full data set, and select the best seven-variable model. It is important that we make use of the full
data set in order to obtain more accurate coefficient estimates. Note that we perform best subset selection on the full data set and select the best
seven-variable model, rather than simply using the variables that were obtained from the training set, because the best seven-variable model on
the full data set may differ from the corresponding model on the training set.
 regfit.best <- regsubsets(Salary ~ ., data = Hitters,</pre>
     nvmax = 19)
 coef(regfit.best, 7)
                     Hits
                                  Walks
                                                   CAtBat
                                                                 CHits
                                                                               CHmRun
 ## (Intercept)
 ## 79.4509472 1.2833513 3.2274264 -0.3752350 1.4957073 1.4420538
 ## DivisionW
                     PutOuts
 ## -129.9866432
                     0.2366813
In fact, we see that the best seven-variable model on the full data set has a different set of variables than the best seven-variable model on the
training set.
We now try to choose among the models of different sizes using cross-validation. This approach is somewhat involved, as 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
 n <- nrow(Hitters)</pre>
 set.seed(1)
 folds <- sample(rep(1:k, length = n))
 cv.errors <- matrix(NA, k, 19,</pre>
     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. Note that in the following code R will automatically use our
predict.regsubsets() function when we call predict() because the best.fit object has class regsubsets.
 for (j in 1:k) {
   best.fit <- regsubsets(Salary ~ .,</pre>
        data = Hitters[folds != j, ],
        nvmax = 19)
   for (i in 1:19) {
     pred <- predict(best.fit, Hitters[folds == j, ], id = i)</pre>
     cv.errors[j, i] <-
          mean((Hitters\$Salary[folds == j] - pred)^2)
This has given us a 10 \times 19 matrix, of which the (j,i)th element corresponds to the test MSE for the jth cross-validation fold for the best i-
variable model. We use the apply() function to average over the columns of this matrix in order to obtain a vector for which the ith element is the
cross-validation error for the i-variable model.
 mean.cv.errors <- apply(cv.errors, 2, mean)</pre>
 mean.cv.errors
                     2
                              3
                                       4
                                                 5
                                                          6
           1
 ## 143439.8 126817.0 134214.2 131782.9 130765.6 120382.9 121443.1 114363.7
                    10
                             11
                                      12
                                                13
                                                         14
                                                                   15
 ## 115163.1 109366.0 112738.5 113616.5 115557.6 115853.3 115630.6 116050.0
          17
 ## 116117.0 116419.3 116299.1
 par(mfrow = c(1, 1))
 plot(mean.cv.errors, type = "b")
      140000
```

15 10 Index We see that cross-validation selects a 10-variable model. We now perform best subset selection on the full data set in order to obtain the 10reg.best <- regsubsets(Salary ~ ., data = Hitters,</pre> nvmax = 19) coef(reg.best, 10) CAtBat CRuns ## (Intercept) Hits Walks AtBat ## 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

130000

120000

mean.cv.errors