

statistics, or other traditional measures of model fit on the training data as evidence of a good model fit in the high-dimensional setting. For instance, as we saw in Figure 6.23, one can easily obtain a model with  $R^2 = 1$  when  $p > n$ . Reporting this fact might mislead others into thinking that a statistically valid and useful model has been obtained, whereas in fact this provides absolutely no evidence of a compelling model. It is important to instead report results on an independent test set, or cross-validation errors. For instance, the MSE or  $R^2$  on an independent test set is a valid measure of model fit, but the MSE on the training set certainly is not.

## 6.5 Lab 1: Subset Selection Methods

### 6.5.1 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(ISLR)
> fix(Hitters)
> 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
> 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)
> 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.

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

`regsubsets()`

```

> library(leaps)
> regfit.full=regsubsets(Salary~.,Hitters)
> summary(regfit.full)
Subset selection object
Call: regsubsets.formula(Salary ~ ., Hitters)
19 Variables (and intercept)
...
1 subsets of each size up to 8
Selection Algorithm: exhaustive

```

		AtBat	Hits	HmRun	Runs	RBI	Walks	Years	CAtBat	CHits
1	( 1 )	" "	" "	" "	" "	" "	" "	" "	" "	" "
2	( 1 )	" "	"*	" "	" "	" "	" "	" "	" "	" "
3	( 1 )	" "	"*	" "	" "	" "	" "	" "	" "	" "
4	( 1 )	" "	"*	" "	" "	" "	" "	" "	" "	" "
5	( 1 )	"*	"*	" "	" "	" "	" "	" "	" "	" "
6	( 1 )	"*	"*	" "	" "	" "	"*	" "	" "	" "
7	( 1 )	" "	"*	" "	" "	" "	"*	" "	"*	"*
8	( 1 )	"*	"*	" "	" "	" "	"*	" "	" "	" "

  

		CHmRun	CRuns	CRBI	CWalks	LeagueN	DivisionW	PutOuts
1	( 1 )	" "	" "	"*	" "	" "	" "	" "
2	( 1 )	" "	" "	"*	" "	" "	" "	" "
3	( 1 )	" "	" "	"*	" "	" "	" "	"*
4	( 1 )	" "	" "	"*	" "	" "	"*	"*
5	( 1 )	" "	" "	"*	" "	" "	"*	"*
6	( 1 )	" "	" "	"*	" "	" "	"*	"*
7	( 1 )	"*	" "	" "	" "	" "	"*	"*
8	( 1 )	"*	"*	" "	"*	" "	"*	"*

  

		Assists	Errors	NewLeagueN
1	( 1 )	" "	" "	" "
2	( 1 )	" "	" "	" "
3	( 1 )	" "	" "	" "
4	( 1 )	" "	" "	" "
5	( 1 )	" "	" "	" "
6	( 1 )	" "	" "	" "
7	( 1 )	" "	" "	" "
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,nvmax=19)
> reg.summary=summary(regfit.full)

```

The `summary()` function also returns  $R^2$ , RSS, adjusted  $R^2$ ,  $C_p$ , and BIC. We can examine these to try to select the *best* overall model.

```

> names(reg.summary)
[1] "which" "rsq" "rss" "adjr2" "cp" "bic"
[7] "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  $R^2$  statistic increases monotonically as more variables are included.

```
> reg.summary$rsq
[1] 0.321 0.425 0.451 0.475 0.491 0.509 0.514 0.529 0.535
[10] 0.540 0.543 0.544 0.544 0.545 0.545 0.546 0.546 0.546
[19] 0.546
```

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="l"` option tells **R** to connect the plotted points with lines.

```
> 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")
```

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  $R^2$  statistic.

`points()`

```
> which.max(reg.summary$adjr2)
[1] 11
> points(11,reg.summary$adjr2[11], col="red",cex=2,pch=20)
```

In a similar fashion we can plot the  $C_p$  and BIC statistics, and indicate the models with the smallest statistic using `which.min()`.

`which.min()`

```
> plot(reg.summary$cp,xlab="Number of Variables",ylab="Cp",
      type='l')
> which.min(reg.summary$cp)
[1] 10
> 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)
```

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 `?plot.regsubsets`.

```
> plot(regfit.full,scale="r2")
> plot(regfit.full,scale="adjr2")
> plot(regfit.full,scale="Cp")
> plot(regfit.full,scale="bic")
```

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
    91.512    -1.869     7.604     3.698     0.643
  DivisionW    PutOuts
 -122.952     0.264
```

### 6.5.2 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,nvmax=19,
  method="forward")
> summary(regfit.fwd)
> regfit.bwd=regsubsets(Salary~.,data=Hitters,nvmax=19,
  method="backward")
> summary(regfit.bwd)
```

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      Walks      CAtBat      CHits
    79.451     1.283     3.227    -0.375     1.496
  CHmRun  DivisionW    PutOuts
    1.442   -129.987     0.237
> coef(regfit.fwd,7)
(Intercept)      AtBat      Hits      Walks      CRBI
   109.787    -1.959     7.450     4.913     0.854
  CWalks  DivisionW    PutOuts
   -0.305   -127.122     0.253
> coef(regfit.bwd,7)
(Intercept)      AtBat      Hits      Walks      CRuns
   105.649    -1.976     6.757     6.056     1.129
  CWalks  DivisionW    PutOuts
   -0.716   -116.169     0.303
```

### 6.5.3 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),rep=TRUE)
> test=(!train)
```

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

```
> regfit.best=regsubsets(Salary~.,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,])
```

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.

`model.  
matrix()`

```
> 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] 220968 169157 178518 163426 168418 171271 162377 157909
[9] 154056 148162 151156 151742 152214 157359 158541 158743
[17] 159973 159860 160106
> which.min(val.errors)
[1] 10
> coef(regfit.best,10)
(Intercept)      AtBat      Hits      Walks      CAtBat
   -80.275    -1.468     7.163     3.643    -0.186
      CHits     CHmRun    CWalks    LeagueN   DivisionW
     1.105     1.384    -0.748    84.558   -53.029
    PutOuts
     0.238
```

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. 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 ten-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 ten-variable model, rather than simply using the variables that were obtained from the training set, because the best ten-variable model on the full data set may differ from the corresponding model on the training set.

```
> regfit.best=regsubsets(Salary~.,data=Hitters,nvmax=19)
> coef(regfit.best,10)
(Intercept)      AtBat      Hits      Walks      CAtBat
   162.535    -2.169     6.918     5.773    -0.130
      CRuns     CRBI    CWalks    DivisionW   PutOuts
     1.408     0.774    -0.831   -112.380     0.297
    Assists
     0.283
```

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. 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
> 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  $j$ th 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  $10 \times 19$  matrix, of which the  $(i, j)$ th element corresponds to the test MSE for the  $i$ th 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  $j$ th element is the cross-validation error for the  $j$ -variable model. **apply()**

```
> mean.cv.errors=apply(cv.errors,2,mean)
> mean.cv.errors
[1] 160093 140197 153117 151159 146841 138303 144346 130208
[9] 129460 125335 125154 128274 133461 133975 131826 131883
[17] 132751 133096 132805
> 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
      135.751       -2.128        6.924        5.620       -0.139
```

CRuns	CRBI	CWalks	LeagueN	DivisionW
1.455	0.785	-0.823	43.112	-111.146
PutOuts	Assists			
0.289	0.269			

## 6.6 Lab 2: Ridge Regression and the Lasso

We will use the `glmnet` package in order to perform ridge regression and the lasso. The main function in this package is `glmnet()`, which can be used to fit ridge regression models, lasso models, and more. This function has slightly different syntax from other model-fitting functions that we have encountered thus far in this book. In particular, we must pass in an `x` matrix as well as a `y` vector, and we do not use the `y ~ x` syntax. We will now perform ridge regression and the lasso in order to predict `Salary` on the `Hitters` data. Before proceeding ensure that the missing values have been removed from the data, as described in Section 6.5.

```
> x=model.matrix(Salary~.,Hitters)[-1]
> y=Hitters$Salary
```

The `model.matrix()` function is particularly useful for creating `x`; not only does it produce a matrix corresponding to the 19 predictors but it also automatically transforms any qualitative variables into dummy variables. The latter property is important because `glmnet()` can only take numerical, quantitative inputs.

### 6.6.1 Ridge Regression

The `glmnet()` function has an `alpha` argument that determines what type of model is fit. If `alpha=0` then a ridge regression model is fit, and if `alpha=1` then a lasso model is fit. We first fit a ridge regression model.

```
> library(glmnet)
> grid=10^seq(10,-2,length=100)
> ridge.mod=glmnet(x,y,alpha=0,lambda=grid)
```

By default the `glmnet()` function performs ridge regression for an automatically selected range of  $\lambda$  values. However, here we have chosen to implement the function over a grid of values ranging from  $\lambda = 10^{10}$  to  $\lambda = 10^{-2}$ , essentially covering the full range of scenarios from the null model containing only the intercept, to the least squares fit. As we will see, we can also compute model fits for a particular value of  $\lambda$  that is not one of the original `grid` values. Note that by default, the `glmnet()` function standardizes the variables so that they are on the same scale. To turn off this default setting, use the argument `standardize=FALSE`.

Associated with each value of  $\lambda$  is a vector of ridge regression coefficients, stored in a matrix that can be accessed by `coef()`. In this case, it is a  $20 \times 100$