

# R Tutorial for Statistical Learning

Qidi Peng

## 12: Subset selection methods

Code	Comments	Results
<pre>#Delete variables having missing data from a data set. install.packages("ISLR");  library (ISLR); fix(Hitters);  is.na(Hitters\$Salary);  sum(is.na(Hitters\$Salary)); sum(1-is.na(Hitters\$Salary)); Hitters =na.omit(Hitters); sum(is.na(Hitters ));</pre>	<p>The functions detecting and deleting missing data can be found in the package "ISLR".</p> <p>Load the package "ISLR".</p> <p>Check the data frame "Hitters", which is baseball players' personal information. We observe that there are missing data in the variable "Salary".</p> <p>To check which data is missing. True is existing; False means missing. R regards True as 1 and False as 0.</p> <p>Number of existing data in "Salary".</p> <p>Number of missing data in "Salary".</p> <p>Extract the players who have not any missing data.</p> <p>Check the extracted data set. There is no longer missing data.</p>	<p>[1] 59</p> <p>[1] 263</p> <p>[1] 0</p>
<pre>#Best subset methods. install.packages("leaps");  library(leaps); help(regsubsets.formula);  regfit.full=regsubsets(Salary~.,Hitters);  summary(regfit.full);</pre>	<p>The best subset methods function can be found in the package "leaps".</p> <p>Load the package.</p> <p>The function regsubsets.formula() performs best subset method.</p> <p>Let Y=Salary, we perform subset selection among all other variables in the data set Hitters.</p> <p>Print the results. String "*" means the corresponding variable is selected by the method. For instance, this output indicates that the best one-variable model contains only "CRBI" and the best two-variable model contains only "Hits" and "CRBI". In the output, Forced in means whether the variable is contained by all models; Forced out means whether the variable is excluded by all models.</p> <p>By default, regsubsets() only reports results up to the</p>	

<pre> regfit.full=regsubsets(Salary~.,data=Hitters,nvmax =19);  reg.summary=summary(regfit.full); Reg.summary;  #model selection.  names(reg.summary);  reg.summary\$cp;  reg.summary\$bic;  reg.summary\$rsq;  reg.summary\$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",type="l"); which.max(reg.summary\$adjr2); points (11, reg.summary\$adjr2[11], col ="red",cex =2, pch =20);  plot(reg.summary\$cp ,xlab =" Number of Variables ",ylab="Cp", type="l"); which.min(reg.summary\$cp); points (10, reg.summary\$cp [10], col ="red",cex =2, pch =20);  plot(reg.summary\$bic ,xlab=" Number of Variables ",ylab=" BIC", type="l"); which.min(reg.summary\$bic); points(6, reg.summary\$bic[6], col ="red",cex=2,pch =20); </pre>	<p>best eight-variable model. But the nvmax option can be used in order to return as many variables as are desired.</p> <p>Output the results for up to best 19-variable model.</p> <p>Show all the statistics (variables) contained in reg.summary.</p> <p>Show the Cp statistic for all 19 models, the one which minimizes the test MSE is the one containing 10 variables.</p> <p>Show the BIC, the one containing 6 variables is the minimizer of BIC.</p> <p>The R2 statistic is increasing. It is not surprising since R2 is a measure for fitting training data. More variables, more accurate.</p> <p>The adjusted R2 chooses the 11-variable model.</p> <p>Unlike R2, It does the job to measure the test error.</p> <p>Open 4 windows of size 2*2. We are going to illustrate the 4 statistics of all the 19 models: RSS, adjusted R2, Cp and BIC.</p> <p>In Window (1,1), present the plot of RSS. Since RSS is to measure training error, it is decreasing as number of variables increases.</p> <p>In Window (1,2), present the adjusted R2, it has a maximum.</p> <p>Find the maximizer of the adjusted R2.</p> <p>Paint this maximum in red.</p> <p>In Window (2,1), present Cp. The best model minimizes Cp.</p> <p>Find the minimizer of Cp.</p> <p>Paint the minimum value point of Cp in red.</p> <p>In Window (2,2), plot BIC VS number of variables.</p> <p>Find the minimizer of BIC.</p> <p>Add this minimum point in red.</p> <p>We see that all these statistics propose different models. It is reasonable because each statistic has a different statistical point of view.</p> <p>Coefficient estimates of the best 6-variable model.</p>	<p>[1] 11</p> <p>[1] 10</p> <p>[1] 6</p>
---	---	--

<pre>coef(regfit.full,6);</pre>		
<pre>#Forward and backward stepwise selection methods. 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); coef(regfit.full,7); coef(regfit.fwd ,7); coef(regfit.bwd,7);</pre>	<p>In regsubsets(), the argument “method=“forward”” can perform forward stepwise selection.</p> <p>Output the results.</p> <p>The argument “method=“backward”” can perform backward stepwise selection</p> <p>Output the results..</p> <p>We compare the best subsets, forward and backward stepwise methods. We observe that they provide the same one to six-variables models, but the seven-variable models are different. This shows the 3 methods are quite accurate when dealing with small number of predictors.</p>	
<pre># Cross validation. fix(predict.regsubsets); #In the separate pad write: 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}  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) ));  for(j in 1:k){best.fit=regsubsets (Salary~.,data=Hitters[folds !=j,], nvmax =19);for(i in 1:19){pred=predict.regsubsets(best.fit,Hitters[folds==j,], id=i); cv.errors [j,i]=mean((Hitters\$Salary[folds ==j]-pred)^2)};  mean.cv.errors =apply(cv.errors ,2, mean); mean.cv.errors;</pre>	<p>There is no predict() function for the subsets methods, we create one here.</p> <p>10-fold cross validation.</p> <p>Compute the test errors for each subset, using best subsets method.</p> <p>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.</p>	

<pre>which.min(mean.cv.errors);</pre>	<p>This function tells you which model does the cross-validation selects. Attention, the result is random.</p>	
---------------------------------------	--	--

## 13: Ridge regression and the lasso

Codes	Comments	Results
<pre># Ridge regression. We perform a regression y~x.  x=model.matrix (Salary~.,Hitters )[, -1];  y=Hitters\$Salary;  install.packages("glmnet"); library (glmnet); grid =10^ seq (10,-2, length =100);  ridge.mod =glmnet(x,y,alpha =0,lambda=grid);  help(glmnet); coef(ridge.mod)[,50];  ridge.mod\$lambda [50]; predict(ridge.mod,s=50,type ="coefficients")[1:20,];  #subset validation.</pre>	<p>Define predictor set (including all predictors, the salaries of players are declined and the qualitative variables are converted to dummy variables).</p> <p>Define the response y.</p> <p>Ridge and the lasso are contained in the package “glmnet”. We install it.</p> <p>We choose a grid for the tuning parameter lambda, from <math>10^{10}</math> to <math>10^{-2}</math>. The mesh is not equal size.</p> <p>Perform a ridge regression for <math>y \sim x</math>. When the argument <math>\alpha=0</math>, we fit ridge regression; when <math>\alpha=1</math>, we fit the lasso regression; when <math>0 &lt; \alpha &lt; 1</math>, we fit a ridge-the lasso mixture model.</p> <p>To check more information on glmnet().</p> <p>The estimates of the coefficients of the 50<sup>th</sup> lambda value (lambda=11498 in this case.)</p> <p>The corresponding lambda.</p> <p>Predict the first 20 estimates of the coefficients for lambda=50.</p>	<p>[1] 11497.57</p>

<pre> set.seed (1);  train=sample (1: nrow(x), nrow(x)/2);  test=(- train );  y.test=y[test];  ridge.mod =glmnet (x[train ],y[train],alpha =0, lambda =grid , thresh =1e-12);  ridge.pred=predict (ridge.mod ,s=4, newx=x[test,]);  mean(( ridge.pred -y.test)^2); </pre>	<p>Set a seed for subset validation.</p> <p>Select a training data set index.</p> <p>Define test data set index.</p> <p>Define test data set.</p> <p>Perform ridge regression using training data set.</p> <p>Predict the regression, using lambda=4 (must be the same length as the test set).</p> <p>Calculate the test MSE. The result is random.</p>	[1] 101036.8
<pre> # The lasso.  lasso.mod =glmnet (x[train ],y[train],alpha =1, lambda =grid);  plot(lasso.mod);  coef(lasso.mod)[,50];  #Cross validation for choosing lambda; subset validation for #calculating the test MSE.  set.seed (1);  help(cv.glmnet);  cv.out =cv.glmnet (x[train ],y[train],alpha =1);  plot(cv.out);  bestlam =cv.out\$lambda.min;  bestlam;  lasso.pred=predict (lasso.mod ,s=bestlam ,newx=x[test ,]);  mean(( lasso.pred -y.test)^2); </pre>	<p>Perform the lasso on the training data (take alpha=1).</p> <p>Illustrate the results. We see most of the coefficients are around 0.</p> <p>Check the coefficients for lambda=11498. Most are zero.</p> <p>Check the cross validation function for glmnet.</p> <p>Run a 10-fold (by default) cross validation for the lasso regression on the training data set, in order to choose lambda.</p> <p>Illustrate the results. We see how the test MSE behaves as log(lambda) increases.</p> <p>Return the best lambda among the grid.</p> <p>Make the lasso prediction using the best lambda.</p> <p>Calculate the test MSE by using subset validation.</p> <p>The result is random.</p>	<p>[1] 32.18284</p> <p>[1] 101918.2</p>

## 14: Principal component regression (PCR)

Codes	Comments	Results
<pre> #Principal component regression.  install.packages("pls");  library(pls); </pre>	<p>Install the package which contains principal component regression function pcr().</p>	

<pre> pcr.fit=pcr(Salary~., data=Hitters ,scale=TRUE ,validation="CV"); validationplot(pcr.fit ,val.type="MSEP");  summary(pcr.fit);  #Calculate the test MSE using subset validation. pcr.pred=predict (pcr.fit ,x[test ,], ncomp =7);  mean((pcr.pred -y.test)^2);  pcr.fit=pcr(y~x,scale =TRUE ,ncomp =7); summary(pcr.fit); </pre>	<p>One directly performs the PCR.</p> <p>Plot the cross validation MSE VS number of components (predictors). Attention, this cross validation method is used to select the principal components.</p> <p>Output the results. From the summary we see that the best model contains 7 predictors.</p> <p>Use the 7 principle components, we make prediction using the training data set.</p> <p>Calculate the test MSE.</p> <p>Perform the PCR only with 7 components.</p> <p>Print the results.</p>	<p>[1] 85199.48</p>
--	---	---------------------

## 15: Tree-Based Methods

Codes	Comments	Results
<pre> #Classification tree.. install.packages("tree");  library(tree); library(ISLR); attach(Carseats ) ; High=ifelse(Carseats\$Sales &lt;=8,"No","Yes ");  Carseats =data.frame(Carseats ,High);  set.seed(2); train=sample (1:nrow(Carseats), 200); Carseats .test=Carseats[-train ,] ; High.test=High[-train]; tree.carseats =tree(High~.-Sales ,Carseats ,subset=train) ; summary(tree.carseats); plot(tree.carseats); </pre>	<p>Install the package which contains tree-based functions for classification and regression.</p> <p>Consider dataset” Carseat”s for classification. ifelse() function is used to create a variable, which takes on a value of Yes if the Sales variable exceeds 8, and takes on a value of No otherwise. Here we artificially create response labels (2~classes) , called High, for each predictor.</p> <p>Create a data matrix of (X,Y), predictors and responses.</p> <p>Select 200 training data.</p> <p>The rest data are for testing. It is validation set.</p> <p>Run classification tree.</p> <p>Output the results.</p> <p>Plot the results.</p> <p>Predict the rest data.</p> <p>Test the prediction by using validation set. Display</p>	

<pre>tree.pred=predict(tree.carseats ,Carseats .test ,type="class") ; table(tree.pred ,High.test)</pre>	<p>the results by confusion matrix.</p>	<p>High.test tree.pred</p> <table><tr><td>No</td><td>Yes</td></tr><tr><td>No</td><td>86 27</td></tr><tr><td>Yes</td><td>30 57</td></tr></table>	No	Yes	No	86 27	Yes	30 57
No	Yes							
No	86 27							
Yes	30 57							
<pre>#Regression tree..  library(MASS) ;  set.seed(1) ;  train = sample (1:nrow(Boston), nrow(Boston)/2) ;  tree.boston=tree(medv~.,Boston , subset=train) ;  summary(tree.boston);  plot(tree.boston);   </pre>								