

Penalized Regression

G. de los Campos

11/25/2024

So far, we have considered methods that estimate parameters by either maximizing the likelihood function (ML) or by minimizing the residual sum of squares (OLS). These methods have good statistical properties (e.g., OLS is unbiased and has minimum variance among the class of linear unbiased estimators, ML is asymptotically unbiased and asymptotically efficient). However, the performance of these methods can be sub-optimal when the number of parameters to be estimated (e.g., the number of regression coefficients) is large relative to sample size and when predictors are highly colinear. One way to circumvent this problem is by using variable selection procedures.

1) Penalized Regression

In a penalized regression, estimates are obtained by minimizing a penalized log-likelihood or, in the case of linear models, a penalized residual sum of squares:

$$\hat{\beta} = \operatorname{argmin}\{(y - X\beta)'(y - X\beta) + \lambda J(\beta)\}$$

where $J(\beta)$ is a penalty function.

Choosing $\lambda = 0$ leads Ordinary Least Squares estimates; however, for any $\lambda > 0$ the solution won't be equivalent to OLS.

Common choices for the penalty function are the

- L2-norm, $J(\beta) = \sum_j \beta_j^2$ (aka Ridge Regression, Hoerl and Kennard 1970),
- L1 norm $J(\beta) = \sum_j |\beta_j|$ (aka Lasso, Tibshirani, 1996), and,
- Linear combinations of the two $J(\beta) = (1 - \alpha) \sum_j \beta_j^2 + \alpha \sum_j |\beta_j|$ (aka Elastic Net, Zhou and Hastie, 2005), for some $\alpha \in [0, 1]$.

Ridge regression shrunk OLS estimates towards zero, without making variable selection. Lasso and Elastic Net combine variable selection and shrinkage.

Commonly, these models are fitted over a grid of values of the regularization parameter (λ) and optimal value for the penalization parameter is often chosen by evaluating the ability of the fitted models to predict data that was not used to train the models (i.e., testing data).

Best subset selection is obtained using as a penalty a function that counts the number of non-zero effects: $J(\beta) = \sum_j 1(\beta_j \neq 0)$. This approach identifies the best subset of predictors and estimate their coefficients without any shrinkage. However, the optimization problem of best-subset selection is non-convex.

Forward regression is an approach where we select predictors sequentially until some stopping criteria is met. This approach can be viewed as an approximation to best-subset selection; however, a forward regression does not guarantee optimal best subset selection.

In this note we will cover forward regression, Ridge Regression, and Lasso.

2) Data

To illustrate these methods we will use a prostate cancer data set which can be found in this link.

```
DATA=read.table('https://hastie.su.domains/ElemStatLearn/datasets/prostate.data',header=TRUE)
head(DATA)
```

```
##          lcavol  lweight age          lbph svi          lcp gleason pgg45          lpsa
## 1 -0.5798185  2.769459  50 -1.386294  0 -1.386294          6      0 -0.4307829
## 2 -0.9942523  3.319626  58 -1.386294  0 -1.386294          6      0 -0.1625189
## 3 -0.5108256  2.691243  74 -1.386294  0 -1.386294          7     20 -0.1625189
## 4 -1.2039728  3.282789  58 -1.386294  0 -1.386294          6      0 -0.1625189
## 5  0.7514161  3.432373  62 -1.386294  0 -1.386294          6      0  0.3715636
## 6 -1.0498221  3.228826  50 -1.386294  0 -1.386294          6      0  0.7654678
##   train
## 1  TRUE
## 2  TRUE
## 3  TRUE
## 4  TRUE
## 5  TRUE
## 6  TRUE
```

```
train=DATA[, 'train']
DATA=DATA[, -ncol(DATA)]
```

In the handout, we will fit various models. To compare them, in the INCLASS assignments you will be asked to estimate the prediction accuracy of each of the models in the testing set. For that evaluation we will use the following partition into training and testing data.

```
“r ## Training and testing data DATA.TRN=DATA[train,] DATA.TST=DATA[!train,]
dim(DATA.TRN) “
## [1] 67  9
r dim(DATA.TST)
## [1] 30  9
```

Our objective is to build a prediction model for the log of the PSA (**lpsa**) as a function of the other variables included in the data set.

3) Full model fitted via OLS

First, we fit a model using all the available predictors via OLS to the training data.

```
fm=lm(lpsa~.,data=DATA.TRN)
summary(fm)
```

```
##
## Call:
## lm(formula = lpsa ~ ., data = DATA.TRN)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.64870 -0.34147 -0.05424  0.44941  1.48675
##
## Coefficients:
```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.429170   1.553588   0.276  0.78334
## lcavol      0.576543   0.107438   5.366 1.47e-06 ***
## lweight     0.614020   0.223216   2.751  0.00792 **
## age        -0.019001   0.013612  -1.396  0.16806
## lbph       0.144848   0.070457   2.056  0.04431 *
## svi        0.737209   0.298555   2.469  0.01651 *
## lcp       -0.206324   0.110516  -1.867  0.06697 .
## gleason    -0.029503   0.201136  -0.147  0.88389
## pgg45      0.009465   0.005447   1.738  0.08755 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7123 on 58 degrees of freedom
## Multiple R-squared:  0.6944, Adjusted R-squared:  0.6522
## F-statistic: 16.47 on 8 and 58 DF,  p-value: 2.042e-12
```

If we were to use pvalues, we would select `lcavol`, `lweight`, and `svi`; however, due to colinearity, we may be leaving out of the model important variables.

Prediction accuracy in the testing set

(See In-class 20)

4) Forward Regression

The `step()` function can be used to implement forward regression (as well as backward elimination and stepwise methods that combine both forward and backward regression).

The following example shows how to fit a forward regression. In argument `object` provides an initial model to the stepwise procedure (forward in this case), the argument `scope` defines the range of models to be explored. The algorithm will evaluate models between `object` and `scope`, `direction` defines whether to perform forward regression, backward elimination, or *both* regression which explores both forward and backward.

```
fm0=lm(lpsa~1,data=DATA)
fullModel='lpsa ~ lcavol+lweight+age+lbph+svi+lcp+gleason+pgg45'
fwd=step(object=fm0,scope=fullModel,direction='forward',data=DATA)
```

```
## Start:  AIC=28.84
## lpsa ~ 1
##
##           Df Sum of Sq    RSS    AIC
## + lcavol   1     69.003  58.915 -44.366
## + svi      1     41.011  86.907 -6.658
## + lcp      1     38.528  89.389 -3.926
## + lweight  1     24.019 103.899 10.665
## + pgg45    1     22.814 105.103 11.783
## + gleason  1     17.416 110.502 16.641
## + lbph     1      4.136 123.782 27.650
## + age      1      3.679 124.239 28.007
## <none>          127.918 28.838
##
## Step:  AIC=-44.37
## lpsa ~ lcavol
##
##           Df Sum of Sq    RSS    AIC
```

```

## + lweight 1 7.1726 51.742 -54.958
## + svi 1 5.2375 53.677 -51.397
## + lbph 1 3.2658 55.649 -47.898
## + pgg45 1 1.6980 57.217 -45.203
## <none> 58.915 -44.366
## + lcp 1 0.6562 58.259 -43.452
## + gleason 1 0.4156 58.499 -43.053
## + age 1 0.0025 58.912 -42.370
##
## Step: AIC=-54.96
## lpsa ~ lcavol + lweight
##
## Df Sum of Sq RSS AIC
## + svi 1 5.1737 46.568 -63.177
## + pgg45 1 1.8158 49.926 -56.424
## <none> 51.742 -54.958
## + lcp 1 0.8187 50.923 -54.506
## + gleason 1 0.7163 51.026 -54.311
## + age 1 0.6456 51.097 -54.176
## + lbph 1 0.4440 51.298 -53.794
##
## Step: AIC=-63.18
## lpsa ~ lcavol + lweight + svi
##
## Df Sum of Sq RSS AIC
## + lbph 1 0.97296 45.595 -63.226
## <none> 46.568 -63.177
## + age 1 0.62301 45.945 -62.484
## + pgg45 1 0.50069 46.068 -62.226
## + gleason 1 0.34449 46.224 -61.898
## + lcp 1 0.06937 46.499 -61.322
##
## Step: AIC=-63.23
## lpsa ~ lcavol + lweight + svi + lbph
##
## Df Sum of Sq RSS AIC
## + age 1 1.15879 44.437 -63.723
## <none> 45.595 -63.226
## + pgg45 1 0.33173 45.264 -61.934
## + gleason 1 0.20691 45.389 -61.667
## + lcp 1 0.10115 45.494 -61.441
##
## Step: AIC=-63.72
## lpsa ~ lcavol + lweight + svi + lbph + age
##
## Df Sum of Sq RSS AIC
## <none> 44.437 -63.723
## + pgg45 1 0.66071 43.776 -63.176
## + gleason 1 0.47674 43.960 -62.769
## + lcp 1 0.13040 44.306 -62.008

```

Prediction accuracy in testing data

(see In-class 20)

Remarks:

- For large-scale screenings the step function can be too slow. As an alternative you can use the `FWD()` function of the `BGDataExt` package.
- The functions `stepAIC()` and `stepBIC()` of the `MASS` package perform forward regression by minimizing AIC and BIC, respectively.

6) Lasso regression using the `glmnet` package

The following example shows how to fit a Lasso regression using the `glmnet` package.

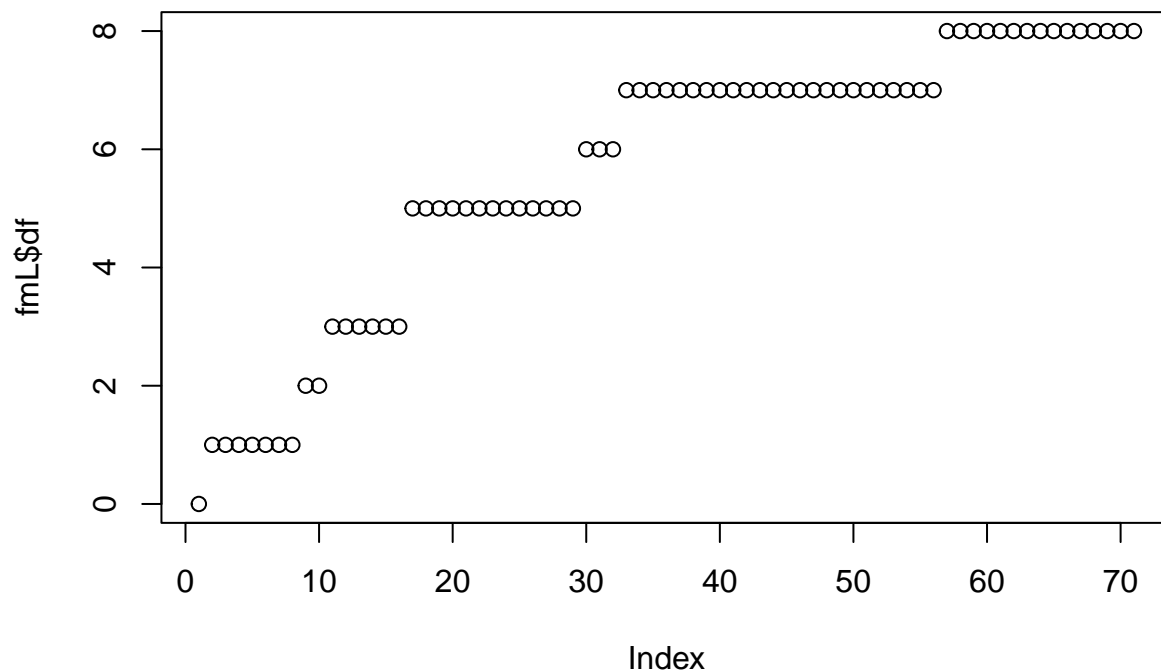
- $\alpha = 1$ is used to fit Lasso (use $\alpha = 0$ for Ridge Regression, and any value between 0 and 1 for Elastic Net).
- By default `glmnet()` runs models over a grid of 100 values of the penalization parameter (λ).
- The object returned by `glmnet()` is a list containing:
 - `$lambda` the grid of values of λ ,
 - `$beta` a matrix with estimated coefficients (rows) for each value of λ in the grid (columns).
 - `$df`, the number of non-zero effects in the solution for each value of λ .
 - `$a0`, the estimated intercept for each value of λ .

```
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## Loaded glmnet 4.1-6
```

```
y=DATA.TRN[, 'lpsa']  
X=as.matrix(DATA.TRN[, colnames(DATA) != 'lpsa'])  
fmL=glmnet(y=y, x=X, alpha=1)  
  
plot(fmL$df)
```



There were changes in the model DF at the following steps

```
which(diff(fmL$df)!=0)
```

```
## [1] 1 8 10 16 29 32 56
```

Finding what predictors were active at each step in the grid

```
for(i in which(diff(fmL$df)!=0)){  
  message('----- step ',i,' -----')  
  print(colnames(X)[fmL$beta[,i]!=0])  
}
```

```
## ----- step 1 -----
```

```
## character(0)
```

```
## ----- step 8 -----
```

```
## [1] "lcavol"
```

```
## ----- step 10 -----
```

```
## [1] "lcavol" "lweight"
```

```
## ----- step 16 -----
```

```
## [1] "lcavol" "lweight" "svi"
```

```
## ----- step 29 -----
```

```
## [1] "lcavol" "lweight" "lbph" "svi" "pgg45"
```

```
## ----- step 32 -----
```

```
## [1] "lcavol" "lweight" "age" "lbph" "svi" "pgg45"
```

```
## ----- step 56 -----
```

```
## [1] "lcavol" "lweight" "age" "lbph" "svi" "lcp" "pgg45"
```

Evaluating prediction accuracy for each value of lambda

(In-class 21)