

Tuesday Precept 6: Penalized Least-Square Regression

Oct 22, 2024

Lecturer: Qishuo Yin

Scribe: Qishuo Yin

1 Ridge Regression and LASSO Regression are Penalized Least-Square Regressions

The key idea of the penalized least-square is to minimize a penalized loss function $Q(\cdot)$ as we have done to the least-square loss function.

- Penalized least-squares:

$$Q(\beta) = \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p p_\lambda(|\beta_j|)$$

where $\beta \in \mathbb{R}^p$ and $p_\lambda(\cdot)$ is the penalized term.

- Ridge regression: L_2 penalty $p_\lambda(\theta) = \lambda|\theta|^2$.

$$Q(\beta) = \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p \beta_j^2$$

- LASSO regression: L_1 penalty $p_\lambda(\theta) = \lambda|\theta|$.

$$Q(\beta) = \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p |\beta_j|$$

2 Code Basis for Ridge Regression and LASSO Regression

Question: How are we going to conduct ridge regression and LASSO regression on a dataset by R codes? Here is a simple example of polynomial regression: Suppose we have a sample `x.random` of size `n = 1000` from the standard normal distribution and `x` is its ordered vector. `y` is a vector also of size `n = 1000` created by adding independent Gaussian variates with mean 0 and standard deviation $\sigma = 2$ to the polynomial regression function

$$\varphi_\beta(x) = x^2 - 10x^3 + x^5$$

```

1 library(glmnet) # package we are going to use for penalized least-square
  regression
2 set.seed(1111)
3 n=1000
4 x.random= array(rnorm(n,mean=0,sd=1), dim=c(n,1))
5 x = x.random[order(x.random)]
6 X = cbind(x,x^2,x^3,x^4,x^5)
7 beta = c(0,1,-10,0,1)
8 Y = X%*%beta + rnorm(n,mean=0,sd=4)

```

Listing 1: set up dataset

2.1 Ridge Regression

For ridge regression, we run the function `glmnet` with the parameter `alpha = 0`. By default, this function runs the ridge regression optimization for a set of 100 values of the hyper-parameter λ . For convenience, we store these values in a vector called `s`.

```

1 ridge = glmnet(X,Y, alpha=0)
2 s=ridge$lambda

```

Listing 2: ridge regression

We plot the scatterplot of `x` and `Y`, and graphs of the regression functions with different values of λ .

```

1 options(repr.plot.width = 8, repr.plot.height = 6)
2
3 # plot scatterplot of x and Y
4 plot(x,Y)
5
6 # plot regression function with lambda = s[1]
7 x.div = seq(from= min(x),to=max(x),length=1000)
8 X.div = cbind(x.div,x.div^2,x.div^3,x.div^4,x.div^5)
9 y.hat = X.div%*% ridge$beta[,1]
10 lines(x.div,y.hat,col=1,lty=1)
11 txt1 = paste("lambda=",signif(s[1],5))
12
13 # plot regression function with lambda = s[33]
14 y.hat = X.div%*% ridge$beta[,33]
15 lines(x.div,y.hat,col=2,lty=1)
16 txt2= paste("lambda=",signif(s[33],5))
17
18 # plot regression function with lambda = s[66]
19 y.hat = X.div%*% ridge$beta[,66]
20 lines(x.div,y.hat,col=3,lty=1)
21 txt3 = paste("lambda=",signif(s[66],5))
22
23 # plot regression function with lambda = s[100]
24 y.hat = X.div%*% ridge$beta[,100]
25 lines(x.div,y.hat,col=4,lty=1)

```

```

26 txt4 = paste("lambda=", signif(s[100], 5))
27
28 # plot regression function with lambda = 0
29 y.hat = predict(ridge, s=0, newx=X.div)
30 lines(x.div, y.hat, col=5, lty=1)
31 txt5 = paste("lambda=", 0.0)
32 legend("bottomleft", c(txt1, txt2, txt3, txt4, txt5), col=c(1:5), lty=c(1, 1, 1, 1, 1))
33 title("Ridge regression for several values of the hyper-parameter lambda")

```

Listing 3: scatterplot and ridge regression function graphs

At the end of this section, we present a way to compute the optimal value of λ as given by 10-fold cross-validation directly from R codes.

```

1 set.seed(1111)
2 cv.ridge = cv.glmnet(X, Y, alpha=0)
3 options(repr.plot.width = 6, repr.plot.height = 4)
4 plot(cv.ridge)
5 cat(paste("Best lambda = ", signif(cv.ridge$lambda.min, 5)))
6 coef(cv.ridge)

```

Listing 4: optimal value of λ in ridge regression

2.2 LASSO Regression

For LASSO regression, we run the function `glmnet` with the parameter `alpha = 0`. And define the different values of λ we want to take by `li`:

```

1 lasso = glmnet(X, Y, alpha=1)
2 s=lasso$lambda
3 ls = length(s)
4 li = c(1, floor(ls/3), floor(2*ls/3), ls)

```

Listing 5: LASSO regression

Then, we plot the scatterplot of x and Y , and graphs of the regression function with different values of λ :

```

1 options(repr.plot.width = 9, repr.plot.height = 6)
2
3 # plot scatterplot of x and Y
4 plot(x, Y)
5
6 # plot regression function with lambda = s[1]
7 x.div = seq(from= min(x), to=max(x), length=1000)
8 X.div = cbind(x.div, x.div^2, x.div^3, x.div^4, x.div^5)
9 y.hat = X.div%*% lasso$beta[, li[1]]
10 lines(x.div, y.hat, col=1, lty=1)
11 txt1 = paste("lambda=", signif(s[li[1]], 5))
12
13 # plot regression function with lambda = s[floor(length(s))/3]
14 y.hat = X.div%*% lasso$beta[, li[2]]

```

```

15 lines(x.div,y.hat,col=2,lty=1)
16 txt2= paste("lambda=",signif(s[li[2]],5))
17
18 # plot regression function with lambda = s[2*floor(length(s))/3]
19 y.hat = X.div%*% lasso$beta[,li[3]]
20 lines(x.div,y.hat,col=3,lty=1)
21 txt3 = paste("lambda=",signif(s[li[3]],5))
22
23 # plot regression function with lambda = s[length(s)]
24 y.hat = X.div%*% lasso$beta[,li[4]]
25 lines(x.div,y.hat,col=4,lty=1)
26 txt4 = paste("lambda=",signif(s[li[4]],5))
27
28 # plot regression function with lambda = 0
29 y.hat = predict(lasso, s=0, newx=X.div)
30 lines(x.div,y.hat,col=5,lty=1)
31 txt5 = paste("lambda=",0.0)
32 legend("bottomleft",c(txt1,txt2,txt3,txt4,txt5),col=c(1:5),lty=c(1,1,1,1,1))
33 title("LASSO regression for several values of the hyper-parameter lambda")

```

Listing 6: scatterplot and LASSO regression function graphs

Similar to ridge regression, there is also a way to compute the optimal value of λ as given by 10-fold cross-validation directly from R codes.

```

1 set.seed(1111)
2 cv.lasso = cv.glmnet(X,Y, alpha=1)
3 cat(paste("Best lambda = ",signif(cv.lasso$lambda.min,5)))
4 coef(cv.lasso)
5 options(repr.plot.width = 6, repr.plot.height = 4)
6 plot(cv.lasso)

```

Listing 7: optimal value of λ in LASSO regression

3 Extensions of Penalized Least-Square Family

- Smoothly clipped absolute deviation(SCAD)

$$p'_\lambda(\theta) = \lambda \left\{ \mathbb{I}(\theta \leq \lambda) + \frac{(a\lambda - \theta)_+}{(a-1)\lambda} \mathbb{I}(\theta > \lambda) \right\}, \quad a > 2$$

with the solution to be

$$\hat{\theta}_{\text{SCAD}}(z) = \begin{cases} \text{sgn}(z)(|z| - \lambda)_+, & \text{when } |z| \leq 2\lambda; \\ \text{sgn}(z)[(a-1)|z| - a\lambda]/(a-2), & \text{when } 2\lambda < |z| \leq a\lambda; \\ z, & \text{when } |z| \geq a\lambda. \end{cases}$$

- Minimax concave penalty (MCP)

$$p'_\lambda(t) = (\lambda - t/a)_+$$

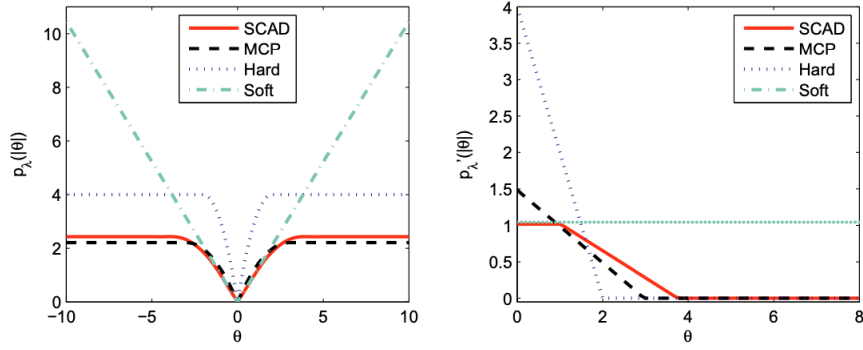


Figure 3.2: Some commonly used penalty functions (left panel) and their derivatives (right panel). They correspond to the risk functions shown in the right panel of Figure 3.3. More precisely, $\lambda = 2$ for hard thresholding penalty, $\lambda = 1.04$ for L_1 -penalty, $\lambda = 1.02$ for SCAD with $a = 3.7$, and $\lambda = 1.49$ for MCP with $a = 2$. Taken from Fan and Lv (2010).

Figure 1: SCAD and MCP

with the solution to be

$$\hat{\theta}_{\text{MCP}}(z) = \begin{cases} \text{sgn}(z)(|z| - \lambda)_+ / (1 - 1/a), & \text{when } |z| < a\lambda; \\ z, & \text{when } |z| \geq a\lambda. \end{cases}$$

The reason why these estimators are introduced may be explained by the following plot from section 3.2 Folded-concave Penalized Least Squares of [1].

References

- [1] J. Fan, R. Li, C.-H. Zhang, and H. Zou. *Statistical foundations of data science*. Chapman and Hall/CRC, 2020.