

Terry crypt.
+ 289

MATH-680

Mohamed 10/10/23

(1)

Regression Shrinkage & Selection via the LASSO.

- The paper focuses on the best subset selection technique using LASSO. ~~and~~ ~~regression shrinkage~~ It also emphasizes on the fact that Lasso minimizes the residual sum of squares in subject to the sum of the absolute value of the coefficient being less than a constant.
- The findings suggest that LASSO enjoys some favourable properties of both subset selection & ridge regression.
- Also the Lasso idea is quite general & can be applied to various statistical models as extensions of generalized linear model or tree-based models.

(2) [2.7]

(a)

For linear regression, the coeff- $\hat{\beta}$ depends on X .

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$$\text{Here } X \text{ is } = \begin{bmatrix} x_1^T \\ x_2^T \\ x_3^T \\ \vdots \\ x_n^T \end{bmatrix}$$

for the model function at $x_0 \rightarrow$

$$\hat{f}(x_0) = x_0^T \hat{\beta} = x_0^T (X^T X)^{-1} X^T y$$

$\left[(X^T X)^{-1} X^T \right]_j$ is the j^{th} row of $(X^T X)^{-1} X^T y$

Then weights are

$$\begin{aligned} l(x_0, X) &= x_{01} \left[(X^T X)^{-1} X^T \right]_1 + x_{02} \left[(X^T X)^{-1} X^T \right]_2 \dots \\ &\quad + x_{0p} \left[(X^T X)^{-1} X^T \right]_p \end{aligned}$$

$$\Rightarrow \sum_{i=1}^p (x_i^T (X^T X)^{-1} x_i^T)$$

So the weights depend on the entire
Set not on y_i .

for K-nearest neighbor reg. function \rightarrow

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in N_k(x)} y_i$$

Weights here are \rightarrow

$$\hat{f}(x) = \sum_i (w_i, x) \geq 1 [i \in N_k(x)]$$

we define $D =$

$$D = \{d(x_i, x_0) : x_i \in T\}$$

Also

$$N_k(x_0) = \{x_i : d(x_i, x_0) \leq d_k\}$$

where d_k is the k th smallest element of D

So for D we need to search the
training set & weights depend on set X .

(b)

~~Total error~~

X_i is fixed and Y varies

No $f(x_0)$ and fixed

$$\text{So, } E y/n (f(x_0) - \hat{f}(x_0))^2$$

$$= f(x_0)^2 - 2 f(x_0) \cdot E y/n (\hat{f}(x_0)) + E y/n ((\hat{f}(x_0))^2)$$

$$= f(x_0) - E y/n (\hat{f}(x_0)) + (E y/n (\hat{f}(x_0)^2) - (E y/n (\hat{f}(x_0)))^2)$$

$$= b(x_0)^2 + \text{var.}(\hat{f}(x_0))$$

(3) [Ex-2.8 page 500]

F or linear regression \rightarrow

$$\text{Training error} = 0.576$$

$$\text{Test error} = 4.121$$

for KNN \rightarrow

$K=1$	Training error = 0	Testing error = 2.473
$K=3$	" " 0.564	" " = 3.022
$K=5$	" " 0.576	" " = 3.022
$K=7$	" " 0.648	" " = 3.022
$K=15$	" " 0.936	" " = 3.846

We can see that KNN has less testing error as compared to linear regression for all the values of K .

Relatively training error is less for $K=1$ & $K=3$ are equal for $K=5$ & linear reg. and more for $K=7$ or $K=15$ than linear reg.

So as the value of K increases the training error increases.

(4) K-NN Classifier

(4)

(a)

K	Turning time	Turning time
1	0	33.4
4	23	32.2
7	27	26.2
10	26	26.4
13	28	25.8
16	30	25.2
30	27.5	21.8
48	28.5	21.9
62	28.5	21.2
80	28.5	21.8
100	28.5	22.3
150	28.5	22.1
200	46	50.1

Plotting R code

(b)	k	Training error	Test Error
	1	0	24.8
	4	18.5	26.2
	7	21	25
	10	23	22.9
	13	22.5	22.2
	16	24.5	22.8
	30	23	23.2
	45	25	23.5
	60	26.5	23.8
	80	28	26.2
	100	27	25.9
	150	30.5	30
	200	53	49.6

Plots in R code

(c)

We choose the k which gives the least testing error. The plots show similar pattern with errors

an increasing training set finally even \rightarrow
the ϵ of K is even. ~~A~~ ~~decides~~ By
the $5/K$ plot the training error is
decreasing

[illegible]

(5)

	k	Training error	Test error
(a)	1	0	2.23
	3	0.42	2.39
	5	0.46	2.39
	7	0.58	2.55
	15	1.04	3.48

(li) LDA training loss = 0.50%

Ln A + estly error = 3.56 %