

# Statistical Machine Learning

## Home Work Three YI CHEN YC3356

1)  $n=2$ ,  $p=2$ ,  $x_{11}=x_{12}=x_1$ ,  $x_{21}=x_{22}=x_2$

(a) 
$$L_2 = \arg \min_{(\hat{\beta}_1, \hat{\beta}_2)} (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + (y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_1)^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

(b) 
$$\frac{\partial L}{\partial \hat{\beta}_1} = \hat{\beta}_1 (x_1^2 + x_2^2 + \lambda) + \hat{\beta}_2 (x_1^2 + x_2^2) - y_1 x_1 - y_2 x_2 = 0$$
  

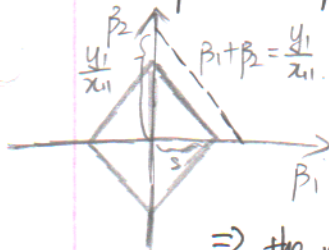
$$\frac{\partial L}{\partial \hat{\beta}_2} = \hat{\beta}_1 (x_1^2 + x_2^2) + \hat{\beta}_2 (x_1^2 + x_2^2 + \lambda) - y_1 x_1 - y_2 x_2 = 0$$

Then we get:  $\hat{\beta}_1 = \hat{\beta}_2$

(c) 
$$L_1 = \arg \min_{(\hat{\beta}_1, \hat{\beta}_2)} (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + (y_2 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + \lambda (|\hat{\beta}_1| + |\hat{\beta}_2|)$$

(d) 
$$\begin{cases} \min (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + (y_2 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 \\ \text{s.t. } |\hat{\beta}_1| + |\hat{\beta}_2| \leq S \end{cases} \Rightarrow \min (y_1 - (\hat{\beta}_1 + \hat{\beta}_2) x_{11})^2$$

This optimization problem has a simple solution  $\hat{\beta}_1 + \hat{\beta}_2 = \frac{y_1}{x_{11}}$



This is a line that parallel to the edge of LASSO diamond  $|\hat{\beta}_1| + |\hat{\beta}_2| = S$ . Now, as we know, the solution must be the contours of function  $(y_1 - (\hat{\beta}_1 + \hat{\beta}_2) x_{11})^2$  that touch the Lasso diamond  $|\hat{\beta}_1| + |\hat{\beta}_2| = S$ .

$\Rightarrow$  the whole line of  $\hat{\beta}_1 + \hat{\beta}_2 = S$  or  $\hat{\beta}_1 + \hat{\beta}_2 = -S$  would be the solution.  
 $\Rightarrow \begin{cases} \hat{\beta}_1 \geq 0, \hat{\beta}_2 \geq 0: \hat{\beta}_1 + \hat{\beta}_2 = S \\ \hat{\beta}_1 \leq 0, \hat{\beta}_2 \leq 0: \hat{\beta}_1 + \hat{\beta}_2 = -S \end{cases} \Rightarrow$  it's not necessary to let  $\hat{\beta}_1 = \hat{\beta}_2$

2) 
$$\hat{g}_1 = \arg \min_g \left( \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int [g^{(3)}(x)]^2 dx \right)$$

$$\hat{g}_2 = \arg \min_g \left( \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int [g^{(4)}(x)]^2 dx \right)$$

(1) As  $\lambda \rightarrow \infty$ ,  $\hat{g}_2$  will have the smaller training RSS.

Reason: the penalty term for  $\hat{g}_2$  is higher in order. Thus, we also need the polynomial of  $g(x)$  to have a higher order to ensure that the deviation exist.

In that way, the model is more flexible in  $\hat{g}_2$  and thus has lower training RSS.

(2) As  $\lambda \rightarrow \infty$ , we already know that the  $\hat{g}_2$  is more flexible and have ~~higher~~ lower RSS. Thus the data has higher probability to overfit. and has higher ~~test~~ test RSS. But, it will also depend on how flexible the data is.

In summary:  $\hat{g}_1$  has smaller test RSS.