

Problem # 1 Ex. 3.12 Show that ridge regression estimates can be obtained by ordinary least squares regression on an augmented data set.

Given

$$X_{new} = \begin{pmatrix} X_{n \times p} \\ \sqrt{\lambda} I \end{pmatrix}_{(n+p) \times (p+1)}$$

$$Y_{new} = \begin{pmatrix} Y_{n \times 1} \\ 0 \end{pmatrix}_{(n+p) \times 1}$$

$$X_{new} \Rightarrow \text{row} = n+p$$

$$\text{column} = p+1$$

$$Y_{new} \Rightarrow \text{row} = n+p$$

$$\text{column} = 1$$

$$\hat{\beta}^{OLS} = (X^T X)^{-1} \cdot X^T \cdot Y$$

we have to show

$$\hat{\beta}^{OLS}_{(X_{new}, Y_{new})} = \hat{\beta}^{ridge}$$

for example  $\Rightarrow$  let taken  $n=3$  and  $p=2$

$$X_{new} = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \\ \sqrt{\lambda} I & \sqrt{\lambda} I \\ \sqrt{\lambda} I & \sqrt{\lambda} I \end{pmatrix} \rightarrow X$$

$$Y_{new} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ 0 \\ 0 \end{pmatrix}$$

$$(X_{\text{new}}^T \cdot X_{\text{new}})$$

$$X_{\text{new}}^T = \begin{pmatrix} X_{11} & X_{21} & X_{31} & \sqrt{\lambda} I & \sqrt{\lambda} I \\ X_{12} & X_{22} & X_{32} & \sqrt{\lambda} I & \sqrt{\lambda} I \end{pmatrix}$$

$\downarrow$   
 $X^T$

$$X_{\text{new}}^T X_{\text{new}} = \begin{pmatrix} X_{11}^2 + X_{21}^2 + X_{31}^2 + 2\lambda I & X_{11}X_{12} + X_{21}X_{22} + X_{31}X_{32} + 2\lambda I \\ X_{12}X_{11} + X_{22}X_{21} + X_{32}X_{31} + 2\lambda I & X_{11}^2 + X_{21}^2 + X_{31}^2 + 2\lambda I \end{pmatrix}$$

In general

$$\left\{ X_{\text{new}}^T X_{\text{new}} = X^T X + 2\lambda I \right\}$$

$P \cdot \lambda = \text{constant}$   
 $(\lambda)$

$$\left[ X_{\text{new}}^T X_{\text{new}} = X^T X + \lambda I \right]$$

$$X_{\text{new}}^T Y_{\text{new}} = \begin{pmatrix} X_{11}Y_1 + X_{21}Y_2 + X_{31}Y_3 + \sqrt{\lambda}I \cdot 0 + \sqrt{\lambda}I \cdot 0 \\ X_{12}Y_1 + X_{22}Y_2 + X_{32}Y_3 + \sqrt{\lambda}I \cdot 0 + \sqrt{\lambda}I \cdot 0 \end{pmatrix}$$

$$X_{\text{new}}^T \cdot Y_{\text{new}} =$$

$$X^T \cdot Y$$

$$0$$

$$\hat{\beta}_{\text{new}}^{\text{OLS}} = (X_{\text{new}}^T \cdot X_{\text{new}}) X_{\text{new}}^T \cdot Y_{\text{new}}$$

$$\left\{ \hat{\beta}_{\text{new}}^{\text{OLS}} = (X^T \cdot X + \lambda I) X^T \cdot Y = \hat{\beta}^{\text{ridge}} \right\}$$

#



## Problem #2 Ex 3.30

Consider the elastic-net optimization problem:

$$\min_{\beta} \|y - X\beta\|^2 + \lambda [\alpha \|\beta\|_2^2 + (1-\alpha) \|\beta\|_1]$$

Show how one can turn this into a lasso problem, using an augmented version of  $X$  and  $y$ .

Solution  $\Rightarrow$

Augmented version of  $X$  and  $y$  will be  $\tilde{X}$  &  $\tilde{y}$

$$\tilde{X} = \begin{bmatrix} X \\ \gamma I_p \end{bmatrix} \quad \tilde{y} = \begin{bmatrix} y \\ 0 \end{bmatrix}$$

$$\tilde{X}\beta = \begin{bmatrix} X\beta \\ \gamma\beta \end{bmatrix}$$

from given hint we know that

$$\left\{ \begin{aligned} \|\tilde{y} - \tilde{X}\beta\|_2^2 &= \left\| \begin{bmatrix} y - X\beta \\ \gamma\beta \end{bmatrix} \right\|_2^2 \\ &= \|y - X\beta\|_2^2 + \gamma^2 \|\beta\|_2^2 \end{aligned} \right\}$$

elastic-net

$$\begin{aligned} \hookrightarrow \min_{\beta} & \left\{ \|y - X\beta\|^2 + \lambda [\alpha \|\beta\|_2^2 + (1-\alpha) \|\beta\|_1] \right\} \\ &= \underbrace{\|y - X\beta\|_2^2}_{\substack{\downarrow \\ \text{Constant} \\ \downarrow \\ \gamma^2}} + \underbrace{\lambda \alpha \|\beta\|_2^2}_{\substack{\downarrow \\ \text{Constant} \\ \downarrow \\ \lambda}} + \lambda (1-\alpha) \|\beta\|_1 \end{aligned}$$

$$= \underbrace{\|Y - X\beta\|_2^2 + \gamma^2 \|\beta\|_2^2}_{\text{}} + \tilde{\lambda} \|\beta\|_1$$

$$= \|\tilde{Y} - \tilde{X}\beta\|_2^2 + \tilde{\lambda} \|\beta\|_1$$

This is a lasso objective function in the form of augment  $\tilde{X}$  and  $\tilde{Y}$



Q4 (3.16) Derive the entries in table 3.4 the explicit forms for estimators in the orthogonal case.

Solution:  $\Rightarrow$

Given table:  $\Rightarrow$

Estimator	Formula.
Best subset (size $M$ )	$\hat{\beta}_j: I( \hat{\beta}_j  \geq  \hat{\beta}_M )$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j) ( \hat{\beta}_j  - \lambda)_+$

by the definition of orthonormal

$$X^T \cdot X = I$$

OLS

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

$$\hat{\beta} = X^T \cdot y$$

① Best subset  $\Rightarrow$  will take the  $m$  predictor with smallest residual sum of square. (RSS).

we know that columns of  $X$  are orthonormal we can construct a basis of Euclidean space  $\mathbb{R}^n$  equipped with the standard inner product. This will be happen by using the first ' $p$ ' columns of  $X$  and the extending there to ' $n-p$ ' linearly independent additional orthonormal vectors.

$$y = \sum_{j=1}^p \hat{\beta}_j x_j + \sum_{j=p+1}^n \gamma_j \tilde{x}_j$$

where  $\hat{\beta}_j$  = Component of  $\hat{\beta}$  in eq

$\gamma_j$  = coefficients of 'j' w.r.t extended basis vector.

Best Subset Selection method estimate of y can be written as

$$\hat{y} = \sum_{j=1}^p I_j \hat{\beta}_j x_j \quad \left\{ \begin{array}{l} \text{where } I_j = 1 \text{ if the} \\ \text{predictor } x_j \text{ are} \\ \text{kept or zero} \\ \text{otherwise.} \end{array} \right.$$

As  $x_i \perp x_j$  are orthonormal  
So,

$$\|y - \hat{y}\|_2^2 = \|y - X\hat{\beta}\|_2^2$$

$$\|y - \hat{y}\|_2^2 = \left\| \underbrace{\sum_{j=1}^p \hat{\beta}_j x_j}_y + \sum_{j=p+1}^N \gamma_j \tilde{x}_j - \underbrace{\sum_{j=1}^p I_j \hat{\beta}_j x_j}_{\hat{y}} \right\|_2^2$$

$$A \|y - \hat{y}\|_2^2 = \left\| \sum_{j=1}^p \hat{\beta}_j x_j (1 - I_j) + \sum_{j=p+1}^N \gamma_j \tilde{x}_j \right\|_2^2$$

$$= \sum_{j=1}^p \hat{\beta}_j^2 (1 - I_j)^2 \|x_j\|_2^2 + \sum_{j=p+1}^N \gamma_j^2 \|\tilde{x}_j\|_2^2$$

$$\|y - \hat{y}\|_2^2 = \sum_{j=1}^p \hat{\beta}_j^2 (1 - I_j)^2 + \sum_{j=1}^N \gamma_j^2$$

We can minimize  $\|y - \hat{y}\|_2^2$  we will choose  $m$  values of  $I_j$  that are equal to one which have the largest values of  $\hat{\beta}_j^2$

Indicator function build a relation between  $A$  and  $x$ .

$1$  = all the elements of  $x$  in  $A$

$0$  = all the elements of  $x$  not in  $A$ .



By the definition of Indicator function ~~our~~  
 we can sort the values of  $|\hat{\beta}_j|$  ~~and~~ and get  
 Only those values with the indices of  
 largest  $m$  meaning where  $I_j=1$  and remaining  
 indices with  $I_j=0$  are taken out.

by using Indicator function

$$\hat{\beta}_j^{\text{best subset}} = \hat{\beta}_j \times I(\text{rank}(|\hat{\beta}_j|) \leq m)$$

$$\left\{ \hat{\beta}_j^{\text{bs}} = \hat{\beta}_j^{\text{ls}} = x^T \cdot y \right\}$$

for Ridge Regression:-

we know that,

$$\hat{\beta}^{\text{ridge}} = \overset{\text{orthonormal}}{(x^T x + \lambda I)^{-1}} x^T \cdot y.$$

$$= (I + \lambda I)^{-1} x^T \cdot y$$

$$= \frac{x^T \cdot y}{1 + \lambda} = \frac{x^T \cdot y}{1 + \lambda} \rightarrow \hat{\beta}^{\text{ls}}$$

$$\boxed{\hat{\beta}^{\text{ridge}} = \frac{\hat{\beta}^{\text{ls}}}{1 + \lambda}}$$

for Lasso:

we know that

$$L(\beta) = (y - x\beta)^T (y - x\beta) + \lambda |\beta|$$

first order derivative w.r.t  $\beta$ .

$$\frac{\partial L(\beta)}{\partial \beta} = -x^T y + x^T x \beta + \lambda \cdot \text{Sign}(\beta)$$

for max  $(\beta)$  will  $\frac{\partial L(\beta)}{\partial \beta} = 0$

$$-x^T y + x^T x \hat{\beta} + \lambda \text{Sign}(\beta) = 0$$

$$\underset{\downarrow}{x^T x} \hat{\beta} = x^T y - \lambda \text{Sign}(\beta)$$

orthonormal

$$I \hat{\beta} = x^T y - \lambda \text{Sign}(\beta)$$

$$\hat{\beta} = I^{-1} (x^T y - \lambda \text{Sign}(\beta)) \quad \therefore I^{-1} = I$$

$$= I (x^T y - \lambda \text{Sign}(\beta))$$

$$\hat{\beta}^{\text{lasso}} = \text{Sign}(\beta) (|x^T y| - \lambda)$$

$$\left\{ \hat{\beta}^{\text{lasso}} = \text{Sign}(\beta) (|x^T y| - \lambda) + \right\}$$