

1. Methods between Ridge and Lasso.

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 .

Let us augment X with $I_{p \times p} * \theta$ to get $\tilde{X} = \begin{bmatrix} X \\ \theta I \end{bmatrix}$

$$\text{Then } \tilde{X}\beta = \begin{bmatrix} X\beta \\ \theta\beta \end{bmatrix}$$

Similarly, let's augment y with p zero values to get
 $\tilde{y} = \begin{bmatrix} y \\ 0 \end{bmatrix}$

$$\implies \|\tilde{y} - \tilde{X}\beta\|_2^2 = \left\| \begin{bmatrix} y - X\beta \\ \theta\beta \end{bmatrix} \right\|_2^2 = \|y - X\beta\|_2^2 + \theta^2\|\beta\|_2^2 \quad (1)$$

Following this augmentation, we estimate β for lasso by

$$\hat{\beta} = \arg \min_{\beta} \left\{ \|\tilde{y} - \tilde{X}\beta\|_2^2 + \tilde{\lambda} \|\beta\|_1 \right\} \quad (2)$$

Combining the two equations, we arrive at

$$\hat{\beta} = \arg \min_{\beta} \left\{ \|y - X\beta\|_2^2 + \theta^2\|\beta\|_2^2 + \tilde{\lambda} \|\beta\|_1 \right\}$$

We transform the above by applying the following transformation:

$$\cdot \theta^2 = \lambda \psi$$

$$\text{Then } \theta = \sqrt{\lambda \psi}$$

$$\cdot \tilde{\lambda} = \lambda(1 - \psi)$$

To optimize the problem, let $\tilde{y} = [y]$, $\tilde{X} = [x]$ and
 $\tilde{\lambda} = \lambda(1 - \psi)$

The lasso problem is solved as such ■

2. Closed forms of Lasso, Ridge, and Best Subset Selection.

Consider a special case of linear regression

$$y = X\beta + \epsilon$$

where the columns of X are orthonormal, that is $X^T X = I$

Let $\hat{\beta}$ denote the least square estimators and $\hat{\beta}_{(M)}$

denote the M^{th} largest absolute value of $\hat{\beta}$, that is

$|\hat{\beta}_{(1)}| \geq |\hat{\beta}_{(2)}| \geq \dots \geq |\hat{\beta}_{(p)}|$. Show the corresponding formulas for Best Subset with size M , Ridge, and Lasso estimator

Estimator	Formula
Best Subset (Size M)	$\hat{\beta}_j * I(\hat{\beta}_j \geq \hat{\beta}_{(M)})$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j) (\hat{\beta}_j - \lambda)_+$

Given that X are orthonormal, $X^T X = I$

For Least Squares, $\hat{\beta} = (X^T X)^{-1} X^T y = X^T y$

(i) Best Subset Selection

First we minimize the residual sum of squares:

$$\Rightarrow (y - X\hat{\beta})^T (y - X\hat{\beta}) = y^T y - 2\hat{\beta}^T X^T y + \hat{\beta}^2$$

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

$$\Rightarrow y^T y - 2(X^T y)^2 + (X^T y)^2 =$$

$$= y^T y - (X^T y)^2$$

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

Thus, the minimal residuals occurs when $\|\hat{\beta}\|^2$ is maximized. The best subset with M predictors adheres to the following property

$$\hat{\beta}_j = I(|\hat{\beta}_j| \geq |\beta_{(M)}|).$$

In real terms this equates to choosing M largest coefficients by their absolute value ■

(ii). Ridge

$$\hat{\beta}_{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \right\}$$

Distributing:

$$= \underset{\beta}{\operatorname{argmin}} \left\{ y^T y - 2y^T X\beta + \cancel{X^T X} \overset{I}{\cancel{\beta^T \beta}} + \lambda \beta^T \beta \right\}$$

$$= \underset{\beta}{\operatorname{argmin}} \left\{ y^T y - 2y^T X\beta + \beta^T \beta (I + \lambda) \right\}$$

Differentiate w.r.t β and set to 0:

$$-2X^T y + 2\beta(I + \lambda) = 0$$

Substitute OLS definition $\hat{\beta}^{\text{OLS}} = X^T y$

$$-2\hat{\beta}^{\text{OLS}} + 2\beta(I + \lambda) = 0$$

$$\beta = \frac{\hat{\beta}^{\text{OLS}}}{I + \lambda} \implies \beta_i = \frac{\hat{\beta}_i}{1 + \lambda}$$

(iii). Closed form of Lasso

$$\hat{\beta}^{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

$$= \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} [y^T y - 2\beta^T X^T y + \cancel{\beta^T X^T X \beta}] + \lambda \|\beta\|_1 \right\}$$

$$= \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} [y^T y - 2\hat{\beta}^{\text{OLS}} + \hat{\beta}^T \hat{\beta}] + \lambda \|\beta\|_1 \right\}$$

Take the derivative w.r.t β and set to 0:

$$-\frac{2\hat{\beta}^{\text{OLS}}}{2} + \frac{2\beta}{2} + \lambda \beta_i = 0$$

We must examine several cases

Case ① $\hat{\beta}_j^{\text{OLS}} \geq 0$ Then $\beta_j \geq 0$

Our equation changes to

$$\begin{aligned} -\hat{\beta}_j^{\text{OLS}} &= \beta_j + \lambda = 0 \\ \hat{\beta}_j^{\text{OLS}} - \lambda &= \beta_j \end{aligned}$$

Note: Our assumption is that $\beta > 0$

$$\text{Then } \beta_j = \text{sign}(\hat{\beta}_j^{\text{OLS}}) (\hat{\beta}_j^{\text{OLS}} - \lambda)_+$$

Case ② $\beta_j < 0$

The expression changes to:

$$-\hat{\beta}_j^{\text{OLS}} + \beta_j - \lambda = 0$$

$$\Rightarrow \hat{\beta}_j^{\text{Lasso}} = \hat{\beta}_j^{\text{OLS}} + \lambda$$

We note that $\hat{\beta}_j^{\text{OLS}} + \lambda$ must be non-negative

Thus we derived the equation:

$$\hat{\beta}_j^{\text{Lasso}} = \text{sign}(\hat{\beta}_j^{\text{OLS}}) * (|\hat{\beta}_j^{\text{OLS}}| - \lambda)_+$$