1. Consider a linear regression model

$$y_i = \sum_{j=1}^{p} x_{ij}\beta_j + \varepsilon_i \quad i = 1, \cdots, n$$

Y and X are centered and standardized, ridge regression is

$$\beta^{ridge} = \arg\min_{\beta}\left(\sum_{i=1}^{n} y_i - \sum_{j=1}^{p} x_{ij}\beta_j\right)^2 + \lambda \sum_{j=1}^{p}\beta_j^2 \quad , \quad \lambda > 0$$

No assumption of X is of full rank. $\beta := [\beta_1, \cdots, \beta_p]^T$

(a) prove that $\beta^{ridge}$ is biased estimator for $\beta$ for given $\lambda$

Denote $P = \|Y - X\beta\|^2 + \lambda\|\beta\|^2$, which is convex, then

$$\frac{\partial P}{\partial \beta} = -2X^T(Y - X\beta) + 2\lambda\beta$$

$$\hat{\beta}^{ridge} = (X^TX + \lambda I)^{-1}X^TY \quad (X^TX + \lambda I) \text{ is always invertible}$$

then $\mathbb{E}\hat{\beta}^{ridge} = \mathbb{E}(X^TX + \lambda I)^{-1}X^TY$

$$= (X^TX + \lambda I)^{-1}X^TX\beta \neq \beta$$

Thus $\hat{\beta}^{ridge}$ is a biased estimator of $\beta$.

(b) find the bias and the variance of $\hat{\beta}^{ridge}$ for given tuning parameter $\lambda$;

$$\text{Bias}(\hat{\beta}^{ridge}) = \mathbb{E}\hat{\beta}^{ridge} - \beta = [(X^TX + \lambda I)^{-1}X^TX - I]\beta$$

$$\text{Var}(\hat{\beta}^{ridge}) = \text{Var}\left[(X^TX + \lambda I)^{-1}X^TY\right]$$

$$= (X^TX + \lambda I)^{-1}X^T[\text{Var }Y]X(X^TX + \lambda I)^{-1}$$

$$= \sigma^2(X^TX + \lambda I)^{-1}X^TX(X^TX + \lambda I)^{-1}$$

(c) Show that $\|\hat{\beta}^{ridge}\|$ increases as the tuning parameter $\lambda \to 0$.

Denote the SVD of X as $X = U\Lambda V^T$ where $\Lambda = \text{diag}(d_1, \cdots, d_p)$, $U \in \mathbb{R}^{n \times p}$ $V \in \mathbb{R}^{p \times p}$

$$U := [u_1, u_2 \cdots u_p] \quad u_i \in \mathbb{R}^n \quad u_i^T u_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

then $X^TX + \lambda I = V\Lambda^2 V^T + V\,diag(\lambda,\cdots,\lambda)V^T$

$$= V\,diag(d_1^2+\lambda,\cdots,d_p^2+\lambda)V^T$$

then $(X^TX+\lambda I)^{-1} = V\,diag\left(\frac{1}{d_1^2+\lambda},\cdots,\frac{1}{d_p^2+\lambda}\right)V^T$

$\hat{\beta}^{ridge} = (X^TX+\lambda I)^{-1}X^TY = V\,diag\left(\frac{1}{d_1^2+\lambda},\cdots,\frac{1}{d_p^2+\lambda}\right)V^TV\Lambda U^TY$

$$= V\,diag\left(\frac{d_1}{d_1^2+\lambda},\cdots\frac{d_p}{d_p^2+\lambda}\right)U^TY$$

$\|\hat{\beta}^{ridge}\|^2 = [\hat{\beta}^{ridge}]^T[\hat{\beta}^{ridge}] = Y^TU\,diag\left(\frac{d_1}{d_1^2+\lambda},\cdots,\frac{d_p}{d_p^2+\lambda}\right)V^TV\,diag\left(\frac{d_1}{d_1^2+\lambda},\cdots,\frac{d_p}{d_p^2+\lambda}\right)U^TY$

$$= Y^TU\,diag\left(\frac{d_1^2}{(d_1^2+\lambda)^2},\cdots\frac{d_p^2}{(d_p^2+\lambda)^2}\right)U^TY.$$

$$= tr\left(Y^TU\,diag\left(\frac{d_1^2}{(d_1^2+\lambda)^2},\cdots,\frac{d_p^2}{(d_p^2+\lambda)^2}\right)U^TY\right)$$

$$= tr\left(diag\left(\frac{d_1^2}{(d_1^2+\lambda)^2},\cdots,\frac{d_p^2}{(d_p^2+\lambda)^2}\right)\underbrace{U^TYY^TU}_{=:\,A(\perp\lambda)}\right)$$

$$= \sum_{i=1}^{p}\frac{d_i^2}{(d_i^2+\lambda)^2}A_{ii}$$

which is a monotonely decreasing function wrt. $\lambda$.

Thus $\|\hat{\beta}^{ridge}\|$ increases as $\lambda \to 0$.


2. Consider the elastic-net optimization problem.

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

(a) Show how the elastic-net optimization problem can turn this into a lasso problem.

Notice the $\|\beta\|_2^2$ can be rewritten as

$$\lambda\alpha\|\beta\|_2^2 = \|\sqrt{\lambda\alpha}\,\beta\|_2^2 = \|0-\sqrt{\lambda\alpha}\,I_p\,\beta\|_2^2$$

We can put it into $\|y-X\beta\|_2^2$ by denoting

$$\tilde{y} = [y^T \; 0\,1_p^T]^T, \quad \tilde{X} = [X^T \; \sqrt{\lambda\alpha}\,I_p]^T$$

then the elastic-net optimization problem can be rewritten as

$$\min_{\beta}\|\tilde{y}-\tilde{X}\beta\|_2^2 + \lambda(1-\alpha)\|\beta\|_1 \quad \text{with the same form of LASSO problem.}$$

(b) Provide your own understanding about the effect of the elastic-net penalty on the param. estimate

From the above alternative problem of elastic-net, we know elastic-net penalty will shrink some of the parameters estimates to 0 just like LASSO penalty by setting $\alpha \neq 1$, and then the other parameters will still have a overall shrinkage encouraged by the augmented data with the same form of ridge penalty.

The turning parameter $\alpha$ is used to adjust the sparsity of parameter estimates.

3. Show the smallest $\lambda$ such that the regression coefficients estimated by the LASSO are all equal to zero is given by $\lambda_{max} = \max_j |\frac{1}{N}<x_j, y>|$

$\hat{\beta}_\lambda = \min_\beta L(\beta) = \min_\beta \frac{1}{2}\|y - X\beta\|^2 + \lambda\|\beta\|_1$ is the objective of LASSO problem.

then $\frac{\partial L}{\partial \beta}|_{\beta = \hat{\beta}_\lambda} = 0 \Rightarrow -X^T(y - X\hat{\beta}_\lambda) + \lambda S(\hat{\beta}_\lambda) = 0$

where $S(\beta_j) = \begin{cases} sgn(\beta_j) & \text{if } \beta_j \neq 0 \\ S_j & \text{if } \beta_j = 0 \end{cases}$ $\quad \exists \; S_j \in [-1, 1]$

Let $\lambda_m$ is a $\lambda$ such that $\hat{\beta}_{\lambda_m} = 0 \in \mathbb{R}^p$

then we have $X^T y - \lambda_m S = 0, \; \exists \; s \in [-1, 1]^p$

Then $\lambda_{max} = \min_{\lambda_m > 0} \lambda_m \quad$ s.t. $\begin{cases} <x_j, y> - \lambda_m S_j = 0, \; j = 1, \ldots, p. \\ -1 \leq S_j \leq 1 \end{cases}$

Since for $j$, $\lambda_m = \frac{1}{S_j}<x_j, y> \geq |<x_j, y>|$

then $\lambda_{max} = \min_{\lambda_m > 0} \{\lambda_m : \lambda_m \geq |<x_j, y>|, \; j = 1, \ldots, p\}$

$= \min_{\lambda_m > 0} \{\lambda_m : \lambda_m \geq \max_j |<x_j, y>|\}$

$= \max_j |<x_j, y>|$

4. $y_i = \sum_{j=1}^{P} x_{ij}\beta_j + r_i + \varepsilon_i$ , $\varepsilon_i \overset{iid}{\sim} N(0,\sigma^2)$, $r = (r_1, \cdots, r_N)$ are unknown constants

Consider minimization of $\min\limits_{\substack{\beta \in \mathbb{R}^P \\ r \in \mathbb{R}^N}} \frac{1}{2}\sum\limits_{i=1}^{N}\left(y_i - \sum\limits_{j=1}^{P}x_{ij}\beta_j - r_i\right)^2 + \lambda\sum\limits_{i=1}^{N}|r_i|$ $\quad\cdots\cdots$ (2)

(a) Show this problem is jointly convex in $\beta$ and $r$

Denote $L(\beta, r) := \frac{1}{2}\sum\limits_{i=1}^{N}\left(y_i - \sum\limits_{j=1}^{P}x_{ij}\beta_j - r_i\right)^2 + \lambda\sum\limits_{i=1}^{N}|r_i|$

Denote $X := \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NP} \end{bmatrix} \in \mathbb{R}^{N\times P}$ , $\beta = [\beta_1\ \beta_2 \cdots \beta_P]^T \in \mathbb{R}^{P\times 1}$

$y = [y_1\ y_2 \cdots y_N]^T \in \mathbb{R}^{N\times 1}$, and we know $r = [r_1\ r_2 \cdots r_N]^{T\,N\times 1}$

then $L(\beta, r) = \frac{1}{2}(y - X\beta - r)^T(y - X\beta - r) + \lambda\|r\|_1$

$\frac{\partial L}{\partial \beta} = -X^T(y - X\beta - r)$ , $\frac{\partial L}{\partial r} = -(y - X\beta - r) + \lambda s(r)$

where $s(r_i) = \begin{cases} sgn(r_i) & r_i \neq 0 \\ s & r_i = 0 \end{cases}$ $\quad \exists s \in [-1,1]$

$\frac{\partial^2 L}{\partial \beta^2} = X^T X \quad \frac{\partial^2 L}{\partial\beta\partial r} = X^T \quad \frac{\partial^2 L}{\partial r\partial\beta} = X \quad \frac{\partial^2 L}{\partial r^2} = I$

So the Hessian matrix is

$$H = \begin{bmatrix} \partial^2 L/\partial\beta^2 & \partial^2 L/\partial\beta\partial r \\ \partial^2 L/\partial r\partial\beta & \partial^2 L/\partial r^2 \end{bmatrix} = \begin{bmatrix} X^T X & X \\ X^T & I \end{bmatrix} \succeq 0$$

thus $L(\beta, r)$ is jointly convex in $\beta$ and $r$.

(2) Huber's Loss function $\rho(t;\lambda) = \begin{cases} \lambda|t| - \lambda^2/2 & \text{if } |t| > \lambda \\ t^2/2 & \text{if } |t| \leq \lambda \end{cases}$

$$\min\limits_{\beta \in \mathbb{R}^P} \sum\limits_{i=1}^{N}\rho\left(y_i - \sum\limits_{j=1}^{P}x_{ij}\beta_j;\lambda\right) \quad\cdots\quad (4)$$

Show that problems (2) & (4) have the same solutions $\hat{\beta}$.

Denote $E(\beta) = \sum_{i=1}^{N} \rho(y_i - \sum_{j=1}^{P} x_{ij}\beta_j; \lambda)$

Since $\dfrac{\partial \rho(y_i - \sum_{j=1}^{P} x_{ij}\beta_j; \lambda)}{\partial \beta_k} = \begin{cases} -\lambda x_{ik} \, \text{sgn}(y_i - \sum_{j=1}^{P} x_{ij}\beta_j) & |y_i - \sum_{j=1}^{P} x_{ij}\beta_j| > \lambda \\ -x_{ik}(y_i - \sum_{j=1}^{P} x_{ij}\beta_j) & |y_i - \sum_{j=1}^{P} x_{ij}\beta_j| \leq \lambda \end{cases}$

$$= -x_{ik}(y_i - \sum_{j=1}^{P} x_{ij}\beta_j - e_i)$$

where $e_i = \begin{cases} y_i - \sum_{j=1}^{P} x_{ij}\beta_j - \lambda \, \text{sgn}(y - \sum_{j=1}^{P} x_{ij}\beta_j) & \text{if } |y_i - \sum_{j=1}^{P} x_{ij}\beta_j| > \lambda \\ 0 & \text{if } |y_i - \sum_{j=1}^{P} x_{ij}\beta_j| \leq \lambda \end{cases}$

then $\dfrac{\partial E}{\partial \beta_k} = \sum_{i=1}^{N} \dfrac{\partial \rho(y_i - \sum_{j=1}^{P} x_{ij}\beta_j; \lambda)}{\partial \beta_k} = -\sum_{i=1}^{N} x_{ik}(y_i - \sum_{j=1}^{P} x_{ij}\beta_j - e_i)$

By (2). $\dfrac{\partial L}{\partial \beta_k} = -\sum_{i=1}^{N} x_{ik}(y_i - \sum_{j=1}^{P} x_{ij}\beta_j - r_i)$

$\dfrac{\partial L}{\partial r_i} = -(y_i - \sum_{j=1}^{P} x_{ij}\beta_j - r_i) + \lambda s(r_i)$

$\Rightarrow (y_i - \sum_{j=1}^{P} x_{ij}\beta_j - r_i) = \lambda s(r_i) \qquad s(r_i) = \begin{cases} \text{sgn}(r_i) & r_i \neq 0 \\ s & r_i = 0 \end{cases} \quad \exists s \in [-1,1].$

if $r_i = 0 \Leftrightarrow -\lambda \leq y_i - \sum_{j=1}^{P} x_{ij}\beta_j \leq \lambda$ i.e. $|y_i - \sum_{j=1}^{P} x_{ij}\beta_j| \leq \lambda$

if $r_i > 0 \Leftrightarrow y_i - \sum_{j=1}^{P} x_{ij}\beta_j > \lambda$

if $r_i < 0 \Leftrightarrow y_i - \sum_{j=1}^{P} x_{ij}\beta_j < -\lambda$ $\Bigg\}$ i.e. $|y_i - \sum_{j=1}^{P} x_{ij}\beta_j| > \lambda$.

it has the same form with the solutions derived from optimization problem (2)