

2.2.2 Theory

$$\beta^{\text{ridge}} = \arg \min_{\beta} \underbrace{\|y - x^T \beta\|^2 + \lambda \|\beta\|_2^2}_{f(\beta)}$$

where $f(\beta) \stackrel{\text{def}}{=} \|y - x^T \beta\|^2 + \lambda \|\beta\|_2^2$

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

$$= y^T y - \underbrace{y^T x^T \beta} - \underbrace{\beta^T x y} + \beta^T x x^T \beta + \lambda \beta^T \beta$$

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

differentiation rule

$$\frac{\partial f}{\partial \beta} = 2 (x x^T + \lambda I) \beta - 2 x y$$

$$\frac{\partial^2 f}{\partial \beta^2} = 2(XX^T + \lambda I)$$

Notice for any vector $u \in \mathbb{R}^p$

$$u^T \frac{\partial^2 f}{\partial \beta^2} u = 2 u^T (XX^T + \lambda I) u$$

$$= 2 u^T X X^T u + 2 \lambda u^T u$$

$$= 2 \|X^T u\|^2 + 2 \lambda \|u\|^2 \geq 0$$

$\frac{\partial^2 f}{\partial \beta^2}$ is pos. semi-def, so f is convex.

It suffices to find local minimizers.

To do this, we solve for the critical

points from $\frac{\partial f}{\partial \beta} = 0$, i.e.

$$(XX^T + \lambda I) \beta^{\text{ridge}} = Xy$$

$$\beta^{\text{ridge}} = (XX^T + \lambda I)^{-1} Xy$$

(In fact,

$$f(\beta) - f(\beta^{\text{ridge}}) = \cancel{\nabla f(\beta^{\text{ridge}})}^{\text{=0}} \cdot (\beta - \beta^{\text{ridge}}) + \frac{1}{2!} (\beta - \beta^{\text{ridge}})^T \underbrace{\frac{\partial^2 f}{\partial \beta^2}}_{\text{pos semi-def}} (\beta - \beta^{\text{ridge}}) \geq 0$$

In conclusion: $\beta^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$
($\lambda > 0$)

Q: Why β^{ridge} reduces over-fitting compared to β^{ls} ?

A: Without Loss of Generality, suppose

$X \in \mathbb{R}^{p \times n}$ has full row rank and $n > p$.

$$\text{Let } X = U \Sigma V^T = U \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \\ & & & 0 \end{pmatrix}_{p \times n} V^T$$

where U, V are orthogonal matrices,
 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > 0$.

For β^{ls} , it satisfies

$$\begin{aligned} XX^T \beta^{ls} &= Xy = X(X^T \beta^* + \epsilon) \\ &= XX^T \beta^* + X\epsilon \end{aligned}$$

$$\begin{aligned} \xRightarrow{\text{SVD}} \underbrace{U \Sigma \Sigma^T U^T}_{\stackrel{\text{def}}{=} \tilde{\beta}^{ls}} \beta^{ls} &= \underbrace{U \Sigma \Sigma^T U^T}_{\stackrel{\text{def}}{=} \tilde{\beta}^*} \beta^* + \underbrace{U \Sigma V^T}_{\stackrel{\text{def}}{=} \tilde{\epsilon}} \epsilon \end{aligned}$$

$$\Rightarrow \Sigma \Sigma^T \tilde{\beta}^{ls} = \Sigma \Sigma^T \tilde{\beta}^* + \Sigma \tilde{\epsilon}$$

$$\text{i.e.} \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_p^2 \end{pmatrix} \tilde{\beta}^{ls} = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_p^2 \end{pmatrix} \tilde{\beta}^* +$$

$$\begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ & & \sigma_p & 0 \end{pmatrix} \tilde{\epsilon}_{p \times n}$$

i.e. for each i , $\sigma_i^2 \tilde{\beta}_i^{ls} = \sigma_i^2 \tilde{\beta}_i^* + \sigma_i \tilde{\epsilon}_i$

$$\Rightarrow \tilde{\beta}_i^{ls} = \tilde{\beta}_i^* + \frac{1}{\sigma_i} \tilde{\epsilon}_i, \quad i=1, \dots, p$$

For β^{ridge} , $(XX^T + \lambda I) \beta^{\text{ridge}} = Xy$
 $= X(X^T \beta^* + \epsilon) = XX^T \beta^* + X\epsilon$

$$\stackrel{\text{SVD}}{\Rightarrow} (U \Sigma \Sigma^T U^T + \lambda I) \beta^{\text{ridge}} =$$

$$U \Sigma \Sigma^T U \beta^* + U \Sigma V^T \epsilon$$

i.e. $\underbrace{U(\Sigma \Sigma^T + \lambda I)U^T}_{\tilde{\beta}^{\text{ridge}}} \beta^{\text{ridge}}$

$$= \cancel{\sqrt{\Sigma}} \Sigma^T \underbrace{U^T \beta^*}_{\tilde{\beta}^*} + \cancel{\sqrt{\Sigma}} \underbrace{V^T \epsilon}_{\tilde{\epsilon}}$$

$$\Rightarrow (\Sigma \Sigma^T + \lambda I) \tilde{\beta}^{\text{ridge}} = \Sigma \Sigma^T \tilde{\beta}^* + \Sigma \tilde{\epsilon}$$

i.e. $\begin{pmatrix} \sigma_1^2 + \lambda & & \\ & \ddots & \\ & & \sigma_p^2 + \lambda \end{pmatrix} \tilde{\beta}^{\text{ridge}}$

$$= \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_p^2 \end{pmatrix} \tilde{\beta}^* + \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \end{pmatrix} \tilde{\epsilon}$$

i.e. for each i , $(\sigma_i^2 + \lambda) \tilde{\beta}_i^{\text{ridge}} = \sigma_i^2 \tilde{\beta}_i^* + \sigma_i \tilde{\epsilon}_i$

$$\Rightarrow \tilde{\beta}_i^{\text{ridge}} = \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \tilde{\beta}_i^* + \frac{\sigma_i}{\sigma_i^2 + \lambda} \tilde{\epsilon}_i$$

$$= \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \tilde{\beta}_i^* + \frac{1}{\sigma_i + \frac{\lambda}{\sigma_i}} \tilde{\epsilon}_i$$

$$= \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \tilde{\beta}_i^* + \frac{1}{(\sqrt{\sigma_i^2} - \sqrt{\frac{\lambda}{\sigma_i^2}})^2 + 2\sqrt{\lambda}} \tilde{\epsilon}_i$$

$$\leq \frac{1}{2\sqrt{\lambda}}$$

Remark: (1) As $\lambda \rightarrow 0$, $\tilde{\beta}^{\text{ridge}} = U^T \beta^{\text{ridge}} \rightarrow \tilde{\beta}^{\text{ls}} = U^T \beta^{\text{ls}}$
 thus $\beta^{\text{ridge}} \rightarrow \beta^{\text{ls}}$.

(2) If $\epsilon = 0$, $\tilde{\beta}^{\text{ridge}} = \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \tilde{\beta}^*$

thus $\beta^{\text{ridge}} = \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \beta^*$.

But $\tilde{\beta}^{\text{ls}} = \tilde{\beta}^*$, thus $\beta^{\text{ls}} = \beta^*$.