

$$y = Bx$$

~~Prob 4~~ Theoretical part TP supervised learning JMA205

$$y = Bx$$

$$(1) E[\tilde{y}] = E[cy] = E[(H+D)y] = E[Hy] + E[Dy]$$

$$E[E[y]] = B \quad E[y] = Bx \quad Bx^{-1}$$

Assuming an unbiased OLS model, we have that:

$$E[Hy] = B$$

we also have that if \tilde{B} is unbiased:

$$E[\tilde{B}] = B = E[Hy] + E[Dy] = B + E[Dy] \Rightarrow E[Dy] = 0$$

$$D B x^{-1} = 0$$

$$Dx = 0 \quad (2)$$

Calculating the variance

$$\text{Var}(cy) = c(\text{Var}(y))c^T = c\sigma^2 c^T = \sigma^2 cc^T$$

↗ scalar

$$C = H + D \text{ and } H = (x^T x)^{-1} x^T y, \text{ so}$$

$$\sigma^2 cc^T = \sigma^2 ((x^T x)^{-1} x^T + D)((x^T x)^{-1} x^T + D)^T$$

$$\sigma^2 cc^T = \sigma^2 (x^T x)^{-1} + \sigma^2 (x^T x)^{-1} (Dx)^T + \sigma^2 (Dx)(x^T x)^{-1} + \sigma^2 DD^T$$

In order for \tilde{B} to be unbiased and hence satisfy $E[Dy] = 0 \Rightarrow E[Dx]B = 0$ with $x^T B = 0$

FROM (2)

$$\text{Var}(\tilde{B}) = \sigma^2 cc^T = \sigma^2 (x^T x)^{-1} + \sigma^2 DD^T$$

$$\text{Var}(\tilde{B}) = \text{Var}(B^*) + \sigma^2 DD^T$$

$$DD^T \geq 0$$

$$\textcircled{2} \quad \beta_{\text{ridge}}^* = (x_c^T x_c + \lambda I)^{-1} x_c^T y_c$$

$$E[\beta_{\text{ridge}}^*] = E[(x_c^T x_c + \lambda I)^{-1} x_c^T y_c]$$

$$E[\beta_{\text{ridge}}^*] = E[(x_c^T x_c + \lambda I)^{-1} x_c^T x_c] \beta$$

The SVD decomposition

$$\begin{aligned} \beta_{\text{ridge}}^* &= (x_c^T x_c + \lambda I)^{-1} x_c^T y_c = ([U D V^T]^T [U D V^T] + \lambda I)^{-1} (U D V^T)^T y_c \\ &\rightarrow \beta_{\text{ridge}}^* = V (D^T D + \lambda I)^{-1} D^T U^T y_c \end{aligned}$$

The variance of a ridge estimator is given by

$$\text{Var}(\beta_{\text{ridge}}^*) = \sigma^2 (x_c^T x_c + \lambda I)^{-1} x_c^T x_c (x_c^T x_c + \lambda I)^{-1}$$

$$\text{If } \lambda > 0 \Rightarrow (x_c^T x_c + \lambda I) > x_c^T x_c$$

$$\Rightarrow (x_c^T x_c + \lambda I)^{-1} < (x_c^T x_c)^{-1}$$

$$\Rightarrow \text{Var}(\beta_{\text{OLS}}) > \text{Var}(\beta_{\text{ridge}}^*)$$

If λ is really close to 0 the model will get similar to a OLS, which means that the bias would be 0

We know that $\beta_{\text{ridge}}^* = (x_c^T x_c + \lambda I)^{-1} x_c^T y_c$. If $x_c^T x_c = I$

$$\beta_{\text{ridge}}^* = (I + \lambda I)^{-1} x_c^T y_c = (1 + \lambda)^{-1} x_c^T y_c$$

$$\beta_{\text{OLS}} = \underbrace{(x_c^T x_c)^{-1}}_{I} x_c^T y_c \Rightarrow \beta_{\text{OLS}} = x_c^T y_c \Rightarrow \beta_{\text{ridge}}^* = \beta_{\text{OLS}} / (1 + \lambda)$$

$$(3) \quad \mathbf{B}_{ELNet}^* = \arg \min_{\mathbf{B}} (\mathbf{y}_c - \mathbf{x}_c \mathbf{B})^T (\mathbf{y}_c - \mathbf{x}_c \mathbf{B}) + \lambda_2 \|\mathbf{B}\|_2^2 + \lambda_1 \|\mathbf{B}\|_1$$

↳ Isn't differentiable in 0, but we know that the subgradient for this func in 0 can be defined as

$$\frac{\partial f}{\partial \mathbf{B}} = \begin{cases} -1, & \text{if } \mathbf{B} < 0 \\ 1, & \text{if } \mathbf{B} > 0 \\ [-1, 1], & \text{if } \mathbf{B} = 0 \end{cases}$$

$$\mathbf{x}_c^T \mathbf{x}_c = \text{Id} \quad \mathbf{B}_{OLS}^* = \mathbf{x}_c^T \mathbf{y}_c$$

$$2\mathbf{B}_{OLS}^* - 2\mathbf{B}(1 - \lambda_2) + \lambda_1 = 0$$

$$\mathbf{B} = \frac{\mathbf{B}_{OLS}^* + \frac{\lambda_1}{2}}{(1 - \lambda_2)}$$