

# ST5227 Assignment 1.

Qi Shuoli (A0274285W) E112403

## Linear Regression

### Problem 1.

Since we let  $X_{\text{new}} = XD$  for some diagonal matrix  $D$

So  $\bar{Y} = XD\beta + \varepsilon$  and loss function  $L(\beta; \bar{Y}, \pi) = (\bar{Y} - XD\beta)^T (\bar{Y} - XD\beta)$

Compute the first and second derivatives:  $\frac{\partial}{\partial \beta} L(\beta; \bar{Y}, XD) = -2DX^T(\bar{Y} - XD\beta) = -2D\bar{X}^T(\bar{Y} - XD\beta)$   
 $(\frac{\partial^2}{\partial \beta \partial \beta^T} = D\bar{X}^T(XD))$

Let  $\frac{\partial}{\partial \beta} L = 0$  so  $-2D\bar{X}^T + 2DX^T X D \beta = 0 \Rightarrow 2DX^T X D \beta = 2D\bar{X}^T \bar{Y}$   
 so  $\hat{\beta}_{\text{new}} = (DX^T X D)^{-1} D\bar{X}^T \bar{Y} = D^{-1} (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{Y} = D^{-1} \hat{\beta}^{\text{OLS}}$

comparing  $\hat{\beta}^{\text{OLS}} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{Y}$ .  $\hat{\beta}_{\text{new}}$  is having another  $D^{-1}$ .

Consequently, standardizing the predictors is not necessary before ordinary least square.

The  $\hat{\beta}_{\text{new}}$  is same as original, but just with  $\hat{\beta}^{\text{OLS}}$  shrink to  $D^{-1}$ . So the final result is not affected. The other property of the result is actually not affected.

### Problem 2.

Let  $\hat{\beta}^{\text{new}} = CY$ . since  $\bar{Y} = X\beta + \varepsilon$  that is  $\hat{\beta}^{\text{new}} = C(X\beta + \varepsilon) = CX\beta + C\varepsilon$

Since  $\hat{\beta}^{\text{new}}$  is an unbiased estimate of  $\beta$ , where  $C$  is a  $p+1 \times n$  matrix

so  $E(\hat{\beta}^{\text{new}}) = \beta = E(CY)$

so  $\hat{\beta}^{\text{new}} X = Y \cdot CX = Y$  so  $CX = I$ ,  $CX$  equal the identity matrix

Consequently,  $\hat{\beta}^{\text{new}} = CX\beta + C\varepsilon = \beta + C\varepsilon$ .

$$\begin{aligned} \text{Var}(\hat{\beta}^{\text{new}}) &= \text{Var}(CY) = \text{Var}(Y)C^T C = \sigma^2 C C^T \\ &= \sigma^2 (C + (X^T X)^{-1} X^T - (X^T X)^{-1} X^T) (C + (X^T X)^{-1} X^T - (X^T X)^{-1} X^T)^T \\ &= \sigma^2 [(X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T [C - (X^T X)^{-1} X^T] + [C - (X^T X)^{-1} X^T] X (X^T X)^{-1} + [C - (X^T X)^{-1} X^T] X (X^T X)^{-1} [C - (X^T X)^{-1} X^T]] \\ &= \sigma^2 [(X^T X)^{-1} + (X^T X)^{-1} (C - (X^T X)^{-1} X^T) X^T + (C - (X^T X)^{-1} X^T) X (X^T X)^{-1} + (C - (X^T X)^{-1} X^T) X (X^T X)^{-1} [C - (X^T X)^{-1} X^T]] \end{aligned}$$

$$\text{Additionally } (C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T) \bar{X} = C \bar{X} - I = 0$$

$$\text{So } \text{Var}(\hat{\beta}^{\text{new}}) = \sigma^2 [(\bar{X}^T \bar{X})^{-1} + (C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)(C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)^T]$$

$$= \text{Var}(\hat{\beta}^{\text{ols}}) + \sigma^2 (C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)(C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)^T$$

Since for any matrix  $A$ , each term of diagonal of  $AA^T$  is greater than or equal to 0.

So  $C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T$  is semi-definite,

$$\text{so } \text{Var}(\hat{\beta}^{\text{new}} - \hat{\beta}^{\text{ols}}) = \sigma^2 (C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)(C - (\bar{X}^T \bar{X})^{-1} \bar{X}^T)^T \geq 0 \Rightarrow \text{Var}(\hat{\beta}^{\text{new}}) \geq \text{Var}(\hat{\beta}^{\text{ols}})$$

Extension to homoskedastic data

Problem 3

(1). The loss function  $L(\beta) = (\bar{Y} - \bar{X}\beta)^T W (\bar{Y} - \bar{X}\beta)$

$\hat{\beta}^{\text{wls}} = \underset{\beta}{\operatorname{argmin}} L(\beta)$  and get the gradient of  $L(\beta)$  with respect to  $\beta$  to equal to 0 to get  $\hat{\beta}^{\text{wls}}$ :

$$\frac{\partial}{\partial \beta} L(\beta) = 0 \Rightarrow -2 \bar{X}^T W (\bar{Y} - \bar{X}\beta) = 0 \Rightarrow \bar{X}^T W \bar{Y} = \bar{X}^T W \bar{X} \beta$$

$$\therefore \hat{\beta}^{\text{wls}} = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \bar{Y}$$

Then,  $\frac{\partial^2}{\partial \beta \partial \beta} L(\beta) = 2 \bar{X}^T W \bar{X} \geq 0$ , so there is the minimum value for  $L(\beta)$  when  $\hat{\beta}^{\text{wls}} = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \bar{Y}$

(2). Bias( $\hat{\beta}^{\text{wls}}$ ) =  $E(\hat{\beta}^{\text{wls}}) - \beta$

$$E(\hat{\beta}^{\text{wls}}) = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W E(\bar{Y}) = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W (\bar{X}\beta + E(\varepsilon))$$

$$E(\varepsilon) = 0 \quad \therefore E(\hat{\beta}^{\text{wls}}) = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \bar{X} \beta = \beta$$

$$\text{so Bias}(\hat{\beta}^{\text{wls}}) = \beta - \beta = 0$$

$$(3). \text{Var}(\hat{\beta}^{\text{wls}}) = \text{Var}[(\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \bar{Y}] = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \text{Var}(\bar{Y}) W \bar{X} (\bar{X}^T W \bar{X})^{-1}$$

Since  $W$  is a diagonal matrix with weights  $w_i = \frac{1}{\sigma^2}$  on diagonal.

Diagonal matrices are symmetric. So  $W = W^T$

$$((\bar{X}^T W \bar{X})^{-1})^T = ((\bar{X}^T W \bar{X})^T)^{-1} = (\bar{X}^T W^T \bar{X}^T)^{-1} = (\bar{X}^T W \bar{X})^{-1}, \text{ so } \bar{X}^T W \bar{X} \text{ also symmetric}$$

$$\text{so } \text{Var}(\hat{\beta}^{\text{wls}}) = (\bar{X}^T W \bar{X})^{-1} \bar{X}^T W \text{Var}(\bar{Y}) W \bar{X} (\bar{X}^T W \bar{X})^{-1}$$

$$\text{For } \text{Var}(T) = \text{Var}(X\beta + \epsilon) = \text{Var}(X\beta) + \text{Var}(\epsilon) + 2\text{Cov}(X\beta, \epsilon)$$

$X\beta$  is not random, error  $\epsilon$  are independent of predictor  $X$ , so

We do have  $\text{Var}(X\beta) = 0$ ,  $\text{Cov}(X\beta, \epsilon) = 0$ . So  $\text{Var}(X\beta) = 0$ ,  $\text{Cov}(X\beta, \epsilon) = 0$ .

So  $\text{Var}(T) = \text{Var}(\epsilon) = n^{-1}$

$$\text{Var}(\hat{\beta}^{\text{ols}}) = (X^T W X)^{-1} X^T W W^T X (X^T W X)^{-1} = (X^T W X)^{-1} X^T W X (X^T W X)^{-1} = (X^T W X)^{-1}$$

Problem 4.

$$\hat{\beta}^{\text{ols}} = (X^T W X)^{-1} X^T W T \quad E(\hat{\beta}^{\text{ols}}) = \beta \quad \text{Var}(\hat{\beta}^{\text{ols}}) = (X^T W X)^{-1}$$

Let  $\hat{\beta}^{\text{new}} = CT$ . We want to show  $\text{Var}(\hat{\beta}^{\text{new}}) \geq \text{Var}(\hat{\beta}^{\text{ols}})$  ( $\beta$  unbiased)

Since  $T = X\beta + \epsilon$ .  $\hat{\beta}^{\text{new}} = C(X\beta + \epsilon) = CX\beta + CE = \beta + CE$

$\text{Var}(\hat{\beta}^{\text{new}}) = \text{Var}(C\beta + \epsilon)$ . since  $\beta$  represent true coefficients that aim to estimate.

Coefficient does not change with the error term  $\epsilon$  or other random variable. So  $\beta$  is considered a constant. So  $\text{Var}(\hat{\beta}^{\text{new}}) = \text{Var}(CE) = C W C^T$

$$\begin{aligned} \text{Var}(\hat{\beta}^{\text{new}}) &= CW C^T = (C + (X^T X)^{-1} X^T - (X^T X)^{-1} X^T) W (C + (X^T X)^{-1} X^T - (X^T X)^{-1} X^T)^T \\ &= n^{-1} \left[ (X^T X)^{-1} X^T X (X^T X)^{-1} + (X^T X)^{-1} X^T [C - (X^T X)^{-1} X^T] + [C - (X^T X)^{-1} X^T] X (X^T X)^{-1} + [C - (X^T X)^{-1} X^T]^T [C - (X^T X)^{-1} X^T] \right] \\ &= n^{-1} \left[ (X^T X)^{-1} + (X^T X)^{-1} (C - (X^T X)^{-1} X^T)^T X^T + (C - (X^T X)^{-1} X^T) X (X^T X)^{-1} + (C - (X^T X)^{-1} X^T) (C - (X^T X)^{-1} X^T)^T \right] \end{aligned}$$

$$\text{Additionally. } n^{-1} (C - (X^T X)^{-1} X^T)^T X^T = C W X - I = CX - I = 0$$

$$\begin{aligned} \text{so } \text{Var}(\hat{\beta}^{\text{new}}) &= (X^T W X)^{-1} + n^{-1} \left[ (C - (X^T X)^{-1} X^T) (C - (X^T X)^{-1} X^T)^T \right] \\ &= \text{Var}(\hat{\beta}^{\text{ols}}) + \left[ n^{-\frac{1}{2}} (C - (X^T X)^{-1} X^T) (C - (X^T X)^{-1} X^T)^T n^{-\frac{1}{2}} \right]^T \end{aligned}$$

Since for any matrix  $A$ , each term of diagonal of  $AA^T$  is greater than or equal to 0.

So  $(C - (X^T X)^{-1} X^T)^T$  is semi-definite,

$$\text{so } \text{Var}(\hat{\beta}^{\text{new}} - \hat{\beta}^{\text{ols}}) = n^{-\frac{1}{2}} (C - (X^T X)^{-1} X^T) (C - (X^T X)^{-1} X^T)^T \geq 0 \Rightarrow \text{Var}(\hat{\beta}^{\text{new}}) \geq \text{Var}(\hat{\beta}^{\text{ols}})$$

Problem 5

$$E[(T - \hat{T})^2] = E[E(T - \hat{f}(x))^2 | \hat{f}(x)] = E[(f(x) - \hat{f}(x))^2] + \text{Var}(\epsilon)$$

$$\begin{aligned}
&= E \left[ f(x) - 2f(x)\hat{f}(x) + \hat{f}^2(x) \right] + \text{var}(\epsilon) \\
&= E(f(x)) - E(2f(x)\hat{f}(x)) + E(\hat{f}^2(x)) + \text{var}(\epsilon) \\
&= E[\hat{f}^2(x)] - 2E[\hat{f}(x)]E[\hat{f}(x)] + \hat{f}^2(x) + \text{var}(\epsilon) \\
&= E[\hat{f}^2(x)] - [E(\hat{f}(x))]^2 + [E(\hat{f}(x))]^2 - 2E[\hat{f}(x)]E[\hat{f}(x)] + \hat{f}^2(x) + \text{var}(\epsilon) \\
&= E[\hat{f}^2(x)] - 2E[\hat{f}(x)]^2 + E[\hat{f}^2(x)]^2 + [E(\hat{f}(x)) - \hat{f}(x)]^2 + \text{var}(\epsilon) \\
&= E[\hat{f}^2(x)] - 2E[\hat{f}(x)]E[\hat{f}(x)] + E[E[\hat{f}(x)]]^2 + [E(\hat{f}(x)) - \hat{f}(x)]^2 + \text{var}(\epsilon) \\
&= [E[\hat{f}(x)] - E(E[\hat{f}(x)])]^2 + [E(\hat{f}(x)) - \hat{f}(x)]^2 + \text{var}(\epsilon) \\
&= E[\hat{f}(x) - E(\hat{f}(x))]^2 + [E(\hat{f}(x)) - \hat{f}(x)]^2 + \text{var}(\epsilon) \\
&= \text{Var}(\hat{f}(x)) + (\text{Bias}(\hat{f}(x)))^2 + \text{Var}(\epsilon)
\end{aligned}$$

Problem b

Let loss function  $\mathcal{L}(\beta) = (Y - X\beta)^T(Y - X\beta) + \lambda\beta^T\beta$

so  $\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \mathcal{L}(\beta)$  and then use gradient to find  $\hat{\beta}^{\text{ridge}}$

$$\frac{\partial}{\partial \beta} (Y - X\beta)^T(Y - X\beta) + \lambda\beta^T\beta = -2X^T(Y - X\beta) + 2\lambda\beta = 0$$

$$\therefore -2X^T(Y - X\beta) + 2\lambda\beta = 0 \Rightarrow X^T(Y - X\beta) = \lambda\beta \Rightarrow X^TY - X^TX\beta = \lambda\beta$$

$$\therefore X^TX\beta + \lambda\beta = X^TY \Rightarrow (X^TX + \lambda I)\beta = X^TY$$

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

$\frac{\partial}{\partial \beta} \mathcal{L}(\beta) = 2X^TX + 2\lambda I \geq 0$  so there is the minimum value for  $\mathcal{L}(\beta)$  when  $\hat{\beta}^{\text{ridge}} = (X^TX + \lambda I)^{-1}X^TY$

Problem 7.

$$\hat{\beta}^{\text{ridge}} = (X^TX + \lambda I)^{-1}X^TY = (X^TX + \lambda I)^{-1}X^T(X\beta + \epsilon)$$

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

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

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

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

$$= (\bar{X}\bar{X} + \lambda I)^{-1} [\bar{X}\bar{X} - (\bar{X}\bar{X} + \lambda I)] \beta = -\lambda (\bar{X}\bar{X} + \lambda I)^{-1} \beta$$

$$\text{Var}(\hat{\beta}^{\text{ridge}}) = (\bar{X}\bar{X} + \lambda I)^{-1} \bar{X} \text{Var}(\epsilon) \bar{X} (\bar{X}\bar{X} + \lambda I)^{-1}$$

$$\text{Since } \text{Var}(\epsilon) = \sigma^2 I, \text{Var}(\hat{\beta}^{\text{ridge}}) = \sigma^2 (\bar{X}\bar{X} + \lambda I)^{-1} \bar{X} \bar{X} (\bar{X}\bar{X} + \lambda I)^{-1}$$

Problem 8.

$$\hat{\beta}^{\text{ols}} = (\bar{X}\bar{X})^{-1} \bar{Y} \quad \hat{\beta}^{\text{ridge}} = (\bar{X}\bar{X} + \lambda I)^{-1} \bar{X} \bar{Y}$$

$$\text{Var}(\hat{\beta}^{\text{ols}}) = (\bar{X}\bar{X})^{-1} \bar{X} \sigma^2 I \bar{X} (\bar{X}\bar{X})^{-1} = (\bar{X}\bar{X})^{-1} \sigma^2$$

$$\text{Var}(\hat{\beta}^{\text{ridge}}) = (\bar{X}\bar{X} + \lambda I)^{-1} \bar{X} \bar{X} (\bar{X}\bar{X} + \lambda I)^{-1} \sigma^2$$

$$\text{Var}(\hat{\beta}) - \text{Var}(\hat{\beta}^{\text{ridge}}) = \sigma^2 (\bar{X}\bar{X} + \lambda I)^{-1} (2\lambda I + \lambda^2 (\bar{X}\bar{X})^{-1}) (\bar{X}\bar{X} + \lambda I)^{-1}$$

$(\bar{X}\bar{X})^{-1}$  is the inverse of a positive definite matrix so  $(\bar{X}\bar{X})^{-1}$  is positive semidefinite.  $y = Ax$   $y^T A^{-1} y = x^T A^T A^{-1} A x = x^T A x \geq 0$

So at the same time  $2\lambda I + \lambda^2 (\bar{X}\bar{X})^{-1}$  should be positive semidefinite.  $\pi^T (A+B)\pi = \pi^T A\pi + \pi^T B\pi \geq 0$

$$\text{Consequently } \text{Var}(\hat{\beta}) - \text{Var}(\hat{\beta}^{\text{ridge}}) = \sigma^2 (\bar{X}\bar{X} + \lambda I)^{-1} (2\lambda I + \lambda^2 (\bar{X}\bar{X})^{-1}) (\bar{X}\bar{X} + \lambda I)^{-1} \geq 0$$

So  $\text{Var}(\hat{\beta}) \geq \text{Var}(\hat{\beta}^{\text{ridge}})$

Problem 9.

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} (\bar{Y} - \bar{X}\beta)^T (\bar{Y} - \bar{X}\beta) + \lambda \beta^T \beta$$

$$X_{\text{new}} = XD$$

$$\hat{\beta}_{\text{new}}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \{ (\bar{Y} - \bar{X}D\beta)^T (\bar{Y} - \bar{X}D\beta) + \lambda \beta^T \beta \}$$

$$\frac{\partial}{\partial \beta} (\bar{Y} - \bar{X}D\beta)^T (\bar{Y} - \bar{X}D\beta) + \lambda \beta^T \beta = -2\bar{X}^T (\bar{Y} - \bar{X}D\beta) + 2\lambda \beta = -2\bar{X}^T (\bar{Y} - \bar{X}D\beta) + 2\lambda \beta = 0$$

$$(XD)^T \bar{Y} - (XD)^T \bar{X} D\beta = \lambda \beta \Rightarrow (XD)^T X D\beta + \lambda \beta = (XD)^T \bar{Y} \Rightarrow ((XD)^T X D + \lambda I)^T \beta = (XD)^T \bar{Y}$$

$$\therefore \hat{\beta}_{\text{new}}^{\text{ridge}} = [(XD)^T X D + \lambda I]^{-1} (XD)^T \bar{Y} = D^{-1} (\bar{X}\bar{X} + \lambda D^{-2})^{-1} \bar{X}^T \bar{Y} \text{ compare } \hat{\beta}^{\text{ridge}} = (\bar{X}\bar{X} + \lambda I)^{-1} \bar{X}^T \bar{Y}$$

Equation be the same if and only  $D^{-2} = I$ .

As ridge regression not linear transform, so we need  $\Rightarrow$  standardize first.