

## THEORETICAL QUESTIONS - ANSWERS

WILLIAM LIAW

### Ordinary Least Squares

Before demonstrating that the Ordinary Least Squares (OLS) estimator has the smallest variance among all linear unbiased estimators, we first analyze the OLS estimator:

$$\begin{aligned}y &= X\beta + \varepsilon \\y - \varepsilon &= X\beta \\X^T(y - \varepsilon) &= X^T X\beta \\(X^T X)^{-1} X^T(y - \varepsilon) &= \beta \\(X^T X)^{-1} X^T y - (X^T X)^{-1} X^T \varepsilon &= \beta \\(X^T X)^{-1} X^T y &= \beta + (X^T X)^{-1} X^T \varepsilon \\\beta^* &= \beta + (X^T X)^{-1} X^T \varepsilon\end{aligned}$$

We can observe that the OLS is unbiased  $\beta^* = \beta$  only if  $X$  is deterministic and  $\mathbb{E}(\varepsilon) = 0$ .

In this case, assuming  $\mathbb{V}(\varepsilon) = \sigma^2 I$ , we can calculate the variance of the OLS estimator:

$$\begin{aligned}\mathbb{V}(\beta^*) &= \mathbb{E}((\beta^* - \mathbb{E}(\beta^*))(\beta^* - \mathbb{E}(\beta^*))^T) \\&= \mathbb{E}((\beta^* - \beta)(\beta^* - \beta)^T) \\&= \mathbb{E}((\beta + (X^T X)^{-1} X^T \varepsilon - \beta)(\beta + (X^T X)^{-1} X^T \varepsilon - \beta)^T) \\&= \mathbb{E}(((X^T X)^{-1} X^T \varepsilon)((X^T X)^{-1} X^T \varepsilon)^T) \\&= \mathbb{E}(((X^T X)^{-1} X^T \varepsilon)(\varepsilon^T X (X^T X)^{-1})) \\&= (X^T X)^{-1} X^T \mathbb{E}(\varepsilon \varepsilon^T) X (X^T X)^{-1} \\&= \sigma^2 (X^T X)^{-1}\end{aligned}$$

Then, we compute the expected value and variance of the alternative estimator  $\tilde{\beta}$ .

### Expected value

First, we compute the expected value of  $\tilde{\beta}$ :

$$\begin{aligned}
\mathbb{E}(\tilde{\beta}) &= \mathbb{E}(Cy) \\
&= C\mathbb{E}(y) \\
&= (H + D)\mathbb{E}(y) \\
&= ((X^T X)^{-1} X^T + D)\mathbb{E}(y) \\
&= ((X^T X)^{-1} X^T + D)\mathbb{E}(X\beta + \varepsilon) \\
&= ((X^T X)^{-1} X^T + D)X\beta \\
&= (I + DX)\beta
\end{aligned}$$

Thus, for  $\tilde{\beta}$  to be unbiased  $(I + DX)\beta = I\beta$ , it is necessary that  $DX = 0$ .

## Variance

Consequently, we compute the variance of  $\tilde{\beta}$ :

$$\begin{aligned}
\mathbb{V}(\tilde{\beta}) &= \mathbb{V}(Cy) \\
&= C\mathbb{V}(y)C^T \\
&= C\mathbb{E}((y - \mathbb{E}(y))(y - \mathbb{E}(y))^T)C^T \\
&= C\mathbb{E}((X\beta + \varepsilon - \mathbb{E}(X\beta + \varepsilon))(X\beta + \varepsilon - \mathbb{E}(X\beta + \varepsilon))^T)C^T \\
&= C\mathbb{E}((\varepsilon - \mathbb{E}(\varepsilon))(\varepsilon - \mathbb{E}(\varepsilon))^T)C^T \\
&= C\mathbb{V}(\varepsilon)C^T \\
&= \sigma^2 C C^T \\
&= \sigma^2 (H + D)(H + D)^T \\
&= \sigma^2 (HH^T + HD^T + DH^T + DD^T) \\
&= \sigma^2 (HH^T + (XX^T)^{-1} X^T D^T + DX(XX^T)^{-1} + DD^T) \\
&= \sigma^2 (HH^T + DD^T) \\
&= \sigma^2 HH^T + \sigma^2 DD^T \\
&= \mathbb{V}(\beta^*) + \sigma^2 DD^T
\end{aligned}$$

## Conclusion

From the last expression, it's evident that  $\mathbb{V}(\tilde{\beta})$  is always greater than or equal to  $\mathbb{V}(\beta^*)$  since  $D$  is non-zero, hence  $DD^T$  introduces additional variance.

The assumption of OLS that we need to use here is that  $X$  is deterministic and  $\mathbb{E}(\varepsilon) = 0$

## Ridge regression

### Biased estimator

To show that the estimator of ridge regression, denoted as  $\beta_{\text{ridge}}^*$ , is biased, we need to demonstrate that its expected value,  $E(\beta_{\text{ridge}}^*)$ , does not equal the true parameter vector  $\beta$ .

Analyzing the ridge estimator  $\beta_{\text{ridge}}^* = \operatorname{argmin}_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$  and denoting  $f(\beta) = \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$ , we have:

$$\begin{aligned} f(\beta) &= \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \\ &= (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \\ &= (y^T - \beta^T X^T)(y - X\beta) + \lambda \beta^T \beta \\ &= (y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta) + \lambda \beta^T \beta \\ \therefore f'(\beta) &= 2X^T X\beta - 2X^T y + 2\lambda\beta \\ \Rightarrow \beta_{\text{ridge}}^* &= (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

Now, let's calculate the expected value of the ridge regression estimator:

$$\begin{aligned} \mathbb{E}(\beta_{\text{ridge}}^*) &= \mathbb{E}((X^T X + \lambda I)^{-1} X^T y) \\ &= \mathbb{E}((X^T X + \lambda I)^{-1} X^T (X\beta + \varepsilon)) \\ &= \mathbb{E}((X^T X + \lambda I)^{-1} X^T X\beta) \\ &= (X^T X + \lambda I)^{-1} X^T X\beta \end{aligned}$$

Therefore, the expected value of the ridge estimator is generally different than  $\beta$  and, consequently, biased. It can be equal to  $\beta$  and unbiased, only if  $\lambda = 0$ , in which case the ridge estimator becomes exactly the OLS estimator.

### SVD decomposition

Given the Singular Value Decomposition (SVD)  $X = UDV^T$ , we can rewrite the expression for  $\beta_{\text{ridge}}^*$ :

$$\begin{aligned}
\beta_{\text{ridge}}^* &= (X^T X + \lambda I)^{-1} X^T y \\
&= ((UDV^T)^T UDV^T + \lambda I)^{-1} (UDV^T)^T y \\
&= ((VDU^T)UDV^T + \lambda I)^{-1} (VDU^T)y \\
&= (VD^2V^T + \lambda I)^{-1} (VDU^T)y \\
&= V(D^2 + \lambda I)^{-1} V^T (VDU^T)y \\
&= V(D^2 + \lambda I)^{-1} DU^T y
\end{aligned}$$

This solution is particularly useful when the matrix  $X$  is ill-conditioned or nearly singular,. In this case, the SVD decomposition provides a numerically stable way to solve the ridge regression problem without directly inverting a potentially singular matrix. Additionally, SVD can be more computationally efficient for large datasets compared to directly computing the inverse of  $X^T X + \lambda I$ , which is no longer necessary through this method.

### Comparison between OLS variance and Ridge variance

First we calculate the variance of the Ridge estimator:

$$\begin{aligned}
\mathbb{V}(\beta_{\text{ridge}}^*) &= \mathbb{V}((X^T X + \lambda I)^{-1} X^T y) \\
&= (X^T X + \lambda I)^{-1} X^T \mathbb{V}(y) (X^T X + \lambda I)^{-1} X^T \\
&= (X^T X + \lambda I)^{-1} X^T \mathbb{V}(y) X (X^T X + \lambda I)^{-1} \\
&= (X^T X + \lambda I)^{-1} X^T \mathbb{V}(\varepsilon) X (X^T X + \lambda I)^{-1} \\
&= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \\
&= \sigma^2 ((UDV^T)^T UDV^T + \lambda I)^{-1} (UDV^T)^T UDV^T ((UDV^T)^T UDV^T + \lambda I)^{-1} \\
&= \sigma^2 ((VDU^T)UDV^T + \lambda I)^{-1} (VDU^T)UDV^T ((VDU^T)UDV^T + \lambda I)^{-1} \\
&= \sigma^2 V(D^2 + \lambda I)^{-1} D^2 (D^2 + \lambda I)^{-1} V^T \\
&= \sum_{i=1}^{\text{rank}(X)} \frac{d_i^2 \sigma^2}{(d_i^2 + \lambda)^2} v_i v_i^T
\end{aligned}$$

We know that the OLS estimator is the Ridge estimator for  $\lambda = 0$ , thus:

$$\mathbb{V}(\beta_{\text{OLS}}^*) = \sum_{i=1}^{\text{rank}(X)} \frac{\sigma^2}{d_i^2} v_i v_i^T$$

Hence, it is apparent that for  $\lambda \geq 0$ , the variance of the Ridge estimator is smaller than or equal to the variance of the OLS estimator, that is  $\mathbb{V}(\beta_{\text{ridge}}^*) \leq \mathbb{V}(\beta_{\text{OLS}}^*)$ .

## Effect of the regularization parameter

Recalling the expression for the expected value of the Ridge estimator:

$$\begin{aligned}\mathbb{E}(\beta_{\text{ridge}}^*) &= (X^T X + \lambda I)^{-1} X^T X \beta \\ &= ((UDV^T)^T UDV^T + \lambda I)^{-1} (UDV^T)^T UDV^T \beta \\ &= ((VDU^T)UDV^T + \lambda I)^{-1} (VDU^T)UDV^T \beta \\ &= (VD^2 V^T + \lambda I)^{-1} VD^2 V^T \beta \\ &= V(D^2 + \lambda I)^{-1} D^2 V^T \beta \\ &= \sum_{i=1}^{\text{rank}(X)} \frac{d_i^2}{d_i^2 + \lambda} v_i v_i^T \beta\end{aligned}$$

One can write the expression for the bias and variance of the ridge estimator:

$$\begin{aligned}b(\beta_{\text{ridge}}^*, \beta) &= \mathbb{E}(\beta_{\text{ridge}}^*) - \beta \\ &= \sum_{i=1}^{\text{rank}(X)} \frac{d_i^2}{d_i^2 + \lambda} v_i v_i^T \beta - \beta \\ \mathbb{V}(\beta_{\text{ridge}}^*) &= \sum_{i=1}^{\text{rank}(X)} \frac{d_i^2 \sigma^2}{(d_i^2 + \lambda)^2} v_i v_i^T\end{aligned}$$

Inference suggests that for small  $\lambda$  values, the ridge estimator closely resembles the OLS estimator, exhibiting low bias but high variance. Conversely, as  $\lambda$  increases, the bias of the ridge estimator amplifies in magnitude while its variance diminishes ( $\lambda \rightarrow +\infty \Rightarrow \mathbb{V}(\beta_{\text{ridge}}^*) \rightarrow 0$ ).

## Derivation of Ridge estimator and OLS estimator expression under $X^T X = I$

As already seen on previous sections:

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

$$\beta_{\text{OLS}}^* = (X^T X)^{-1} X^T y$$

Assuming  $X^T X = I$ , these last expressions become:

$$\begin{aligned}
\beta_{\text{ridge}}^* &= (X^T X + \lambda I)^{-1} X^T y \\
&= (I + \lambda I)^{-1} X^T y \\
&= ((1 + \lambda)I)^{-1} X^T y
\end{aligned}$$

$$\begin{aligned}
\beta_{\text{OLS}}^* &= (X^T X)^{-1} X^T y \\
&= X^T y
\end{aligned}$$

Therefore:

$$\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{1 + \lambda}$$

## Elastic Net

Analyzing the Elastic Net estimator  $\beta_{\text{EINet}}^* = \text{argmin}_{\beta} \|y - X\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$  and denoting  $f(\beta) = \|y - X\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$ , we have:

$$\begin{aligned}
f(\beta) &= \|y - X\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \\
&= (y - X\beta)^T (y - X\beta) + \lambda_2 \beta^T \beta + \lambda_1 \|\beta\|_1 \\
&= (y^T - \beta^T X^T)(y - X\beta) + \lambda_2 \beta^T \beta + \lambda_1 \|\beta\|_1 \\
&= (y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta) + \lambda_2 \beta^T \beta + \lambda_1 \|\beta\|_1 \\
\therefore \partial f(\beta) &= 2X^T X\beta - 2X^T y + 2\lambda_2 \beta \pm \lambda_1 \\
\Rightarrow \beta_{\text{EINet}}^* &= (X^T X + \lambda_2 I)^{-1} (X^T y \mp \frac{\lambda_1}{2})
\end{aligned}$$

Assuming  $X^T X = I$ , this last expression becomes:

$$\begin{aligned}
\beta_{\text{EINet}}^* &= (X^T X + \lambda_2 I)^{-1} X^T (y \mp \frac{\lambda_1}{2}) \\
&= (I + \lambda_2 I)^{-1} X^T (y \mp \frac{\lambda_1}{2}) \\
&= ((1 + \lambda_2)I)^{-1} X^T (y \mp \frac{\lambda_1}{2})
\end{aligned}$$

Therefore:

$$\beta_{\text{EINet}}^* = \frac{\beta_{\text{OLS}}^* \mp \frac{\lambda_1}{2}}{1 + \lambda_2}$$