

# Introduction Supervised Learning

Theoretical questions

## OLS (Ordinary Least Squares)

We have seen that the OLS estimator is equal to  $\beta^* = (X^T X)^{-1} X^T y$  which can be rewritten as  $\beta^* = Hy$ . Let  $\hat{\beta} = Cy$  be another linear unbiased estimator of  $\beta$  where  $C$  is a  $d \times n$  matrix, e.g.,  $C = H + D$  where  $D$  is a non-zero matrix.

- Demonstrate that OLS is the estimator with the smallest variance: compute  $E[\hat{\beta}]$  and  $Var(\hat{\beta}) = E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])^T]$  and show when and why  $Var(\beta^*) < Var(\hat{\beta})$ . Which assumption of OLS do we need to use?

### Answer

To demonstrate that the OLS estimator has the smallest variance, we need to use the Gauss-Markov Theorem, which states that under the assumptions of the classical linear regression model, the OLS estimator has the smallest variance among all unbiased linear estimators. These assumptions are:

1. Linearity of parameters
2. Random sampling
3. No perfect multicollinearity
4. Zero conditional mean (The error term has a zero conditional mean given any value of the explanatory variables)
5. Homoscedasticity (constant variance) of the errors

Given that  $\beta^* = Hy$  and  $\hat{\beta} = Cy$ , where  $C = H + D$ , and assuming that  $H$  is the matrix which gives us the OLS estimator, we have that  $H = (X^T X)^{-1} X^T$ .

For  $\beta^*$ :

$$E[\beta^*] = E[Hy] = HE[y] = HX\beta$$

Since  $H = (X^T X)^{-1} X^T$ , we have  $HX = I$ , where  $I$  is the identity matrix, so  $E[\beta^*] = \beta$ .

For  $\hat{\beta}$ :

$$E[\hat{\beta}] = E[Cy] = CE[y] = CX\beta$$

To be an unbiased estimator,  $E[\hat{\beta}]$  must equal  $\beta$ , which implies that  $CX = I$ .

For the variance of  $\hat{\beta}$ :

$$\begin{aligned} Var(\hat{\beta}) &= E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])^T] \\ Var(\hat{\beta}) &= E[(Cy - CX\beta)(Cy - CX\beta)^T] \\ Var(\hat{\beta}) &= CE[(y - X\beta)(y - X\beta)^T]C^T \end{aligned}$$

Since  $E[(y - X\beta)(y - X\beta)^T]$  is the variance of  $y$ , which we can denote as  $\sigma^2 I$  under the assumption of homoscedasticity and independence, we have:

$$Var(\hat{\beta}) = \sigma^2 CC^T$$

For the variance of the OLS estimator:

$$Var(\beta^*) = \sigma^2 HH^T$$

And since  $H = (X^T X)^{-1} X^T$ , we have:

$$Var(\beta^*) = \sigma^2 (X^T X)^{-1}$$

We need to show that  $Var(\beta^*) < Var(\hat{\beta})$ . Since  $C = H + D$ , we have:

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \sigma^2(H + D)(H + D)^T \\ \text{Var}(\hat{\beta}) &= \sigma^2(HH^T + HD^T + DH^T + DD^T) \end{aligned}$$

Given that  $\text{Var}(\beta^*) = \sigma^2 HH^T$ ,  $\text{Var}(\hat{\beta}) > \text{Var}(\beta^*)$  because  $HD^T + DH^T + DD^T$  is a positive semi-definite matrix, and adding this to  $HH^T$  will give a matrix with larger diagonal elements (variances), assuming  $D$  is not a matrix of zeros (which would violate the assumption that  $D$  is a non-zero matrix). This shows that the variance of  $\beta^*$  is less than the variance of  $\hat{\beta}$ , making  $\beta^*$  the estimator with the smallest variance among all linear unbiased estimators.

## Ridge Regression

Suppose that both  $y$  and the columns of  $x$  are centered ( $y_c$  and  $x_c$ ) so that we do not need the intercept  $\beta_0$ . In this case, the matrix  $x_c$  has  $d$  (rather than  $d + 1$ ) columns. We can thus write the criterion for ridge regression as:

$$\beta_{\text{ridge}}^* = \arg \min_{\beta} \{ (y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda \|\beta\|^2 \}$$

- Show that the estimator of ridge regression is biased (that is  $E[\beta_{\text{ridge}}^*] \neq \beta$ ).

**Answer:** The ridge regression estimator  $\beta_{\text{ridge}}^*$  is found by minimizing the penalized residual sum of squares:

$$\beta_{\text{ridge}}^* = \arg \min_{\beta} \{ (y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda \|\beta\|^2 \}$$

The solution to this minimization problem is:

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

The expectation of  $\beta_{\text{ridge}}^*$ :

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

Since  $y_c = x_c \beta + \epsilon$ , where  $\epsilon$  is the error term, we can substitute  $y_c$  into the expectation:

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

Distributing  $x_c^T$  we get:

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

Assuming  $E[\epsilon] = 0$ , this simplifies to:

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

$E[\beta_{\text{ridge}}^*]$  will not equal  $\beta$  unless  $\lambda = 0$ , because the presence of  $\lambda I$  in the inverse term modifies the relation between  $x_c^T x_c$  and  $\beta$ . Specifically, when  $\lambda > 0$ , the term  $(x_c^T x_c + \lambda I)^{-1} x_c^T x_c$  acts as a shrinkage operator, pulling the estimates of  $\beta$  towards zero.

Therefore, the estimator  $\beta_{\text{ridge}}^*$  is biased because the expectation of the estimator does not equal the true parameter value, i.e.,  $E[\beta_{\text{ridge}}^*] \neq \beta$ .

- Recall that the SVD decomposition is  $x_c = UDV^T$ . Write down by hand the solution  $\beta_{\text{ridge}}^*$  using the SVD decomposition. When is it useful using this decomposition? Hint: do you need to invert a matrix?

**Answer:** Substituting the SVD of  $x_c$  into the expression for  $\beta_{\text{ridge}}^*$  we get:

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

Since  $U^TU = I$  and  $VV^T = I$ , where  $I$  is the identity matrix, we can simplify this to:

$$\beta_{\text{ridge}}^* = (VD^2V^T + \lambda I)^{-1}VDU^T y_c$$

We can take advantage of the diagonal structure of  $D^2$  and the orthogonal matrices  $U$  and  $V$  to compute the ridge estimator more efficiently:

$$\beta_{\text{ridge}}^* = V(D^2 + \lambda I)^{-1}DV^T y_c$$

This is possible because the inverse of a diagonal matrix  $D^2 + \lambda I$  is easy to compute; it's simply the reciprocal of the diagonal elements.

Using the SVD decomposition is particularly useful in ridge regression for a couple of reasons:

1. Numerical stability: When  $x_c^T x_c$  is close to singular or ill-conditioned (which can happen when multicollinearity is present or when  $d$  is large), directly computing its inverse as required in the standard ridge regression formula can be numerically unstable. The SVD approach avoids this problem because the inverse of a diagonal matrix (with the regularization term added) is always well-conditioned.
  2. Computational efficiency: Computing the inverse of a matrix is computationally expensive and can be slow if the matrix is large. However, because SVD provides us with matrices  $U$ ,  $D$ , and  $V$ , where  $D$  is diagonal, we only need to compute the inverse of the diagonal elements of  $D^2 + \lambda I$ , which is straightforward and fast.
- Remember that  $\text{Var}(\beta_{\text{OLS}}^*) = \sigma^2(X^T X)^{-1}$ . Show that  $\text{Var}(\beta_{\text{OLS}}^*) \geq \text{Var}(\beta_{\text{ridge}}^*)$ .

**Answer:** The variance of the OLS estimator is:

$$\text{Var}(\beta_{\text{OLS}}^*) = \sigma^2(X^T X)^{-1}$$

For the ridge regression estimator, the solution can be written using the SVD as  $\beta_{\text{ridge}}^* = V(D^2 + \lambda I)^{-1}DV^T y$ . The variance of the ridge regression estimator is:

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

Now, we need to show that:

$$\sigma^2(X^T X)^{-1} \geq \sigma^2 V(D^2 + \lambda I)^{-2} V^T$$

Using the SVD of  $X$ , we have  $X = UDV^T$ , so  $X^T X = VD^2V^T$ . Replacing this into the variance of the OLS estimator gives us:

$$\text{Var}(\beta_{\text{OLS}}^*) = \sigma^2 (VD^2V^T)^{-1}$$

Multiplying both sides by  $VD^2V^T$  to remove the inverse, we get:

$$VD^2V^T \cdot \text{Var}(\beta_{\text{OLS}}^*) = \sigma^2 I$$

Since  $VD^2V^T$  is a positive semi-definite matrix,  $\text{Var}(\beta_{\text{OLS}}^*)$  must also be a positive semi-definite matrix. This implies that:

$$VD^2V^T \cdot \text{Var}(\beta_{\text{OLS}}^*) \geq \sigma^2 I$$

Similarly, for ridge regression, we have:

$$V(D^2 + \lambda I)^{-2} \cdot \text{Var}(\beta_{\text{ridge}}^*) = \sigma^2 I$$

Multiplying both sides by  $(D^2 + \lambda I)^2$  we get:

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

Given that  $(D^2 + \lambda I)$  is a diagonal matrix with each diagonal element  $d_i^2 + \lambda$  being greater than  $d_i^2$ , the inverse of  $(D^2 + \lambda I)$  will have diagonal elements less than or equal to the inverse of  $D^2$ . Thus:

$$V(D^2 + \lambda I)^{-2} V^T \leq VD^{-2}V^T$$

Multiplying through by  $\sigma^2$  we find:

$$\sigma^2 V(D^2 + \lambda I)^{-2} V^T \leq \sigma^2 VD^{-2}V^T$$

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

Therefore, the variance of the OLS estimator is greater than or equal to the variance of the ridge regression estimator.

- When  $\lambda$  increases what happens to the bias and to the variance? Hint: Compute  $\text{MSE} = E[(y_0 - x_0^T \beta_{\text{ridge}}^*)^2]$  at the test point  $(x_0, y_0)$  with  $y_0 = x_0^T \beta + \epsilon_0$  being the true model and  $\beta_{\text{ridge}}^*$  the ridge estimate.

**Answer:** To examine what happens to the bias and variance as  $\lambda$  increases, let's consider the mean squared error (MSE) at the test point  $(x_0, y_0)$ . The MSE can be decomposed into the sum of the variance and the square of the bias, plus the variance of the error term:

$$\text{MSE} = \text{Var}(x_0^T \beta_{\text{ridge}}^*) + [\text{Bias}(x_0^T \beta_{\text{ridge}}^*)]^2 + \text{Var}(\epsilon_0)$$

Given that  $y_0 = x_0^T \beta + \epsilon_0$ , where  $x_0$  is a new observation and  $\epsilon_0$  is the error term associated with the new observation, the bias of the ridge estimate at this test point is:

$$\text{Bias}(x_0^T \beta_{\text{ridge}}^*) = E[x_0^T \beta_{\text{ridge}}^*] - x_0^T \beta$$

As  $\lambda$  increases, the ridge estimator  $\beta_{\text{ridge}}^*$  will shrink towards zero. This increases the bias term  $E[x_0^T \beta_{\text{ridge}}^*] - x_0^T \beta$  since the expected value of  $x_0^T \beta_{\text{ridge}}^*$  will be further from  $x_0^T \beta$ .

Regarding variance, the ridge estimate's variance is given by:

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

As  $\lambda$  increases, the diagonal elements of the matrix  $(D^2 + \lambda I)$  increase, which leads to a decrease in the diagonal elements of the inverse matrix  $(D^2 + \lambda I)^{-2}$ . Consequently, the variance  $\text{Var}(x_0^T \beta_{\text{ridge}}^*)$  decreases.

As  $\lambda$  increases: - The bias  $\text{Bias}(x_0^T \beta_{\text{ridge}}^*)$  increases because the ridge regression estimate is shrunk further towards zero, causing it to deviate more from the true parameter  $\beta$ . - The variance  $\text{Var}(x_0^T \beta_{\text{ridge}}^*)$  decreases because the regularization term  $\lambda$  penalizes the magnitude of the coefficients, thus reducing the estimator's sensitivity to fluctuations in the training data.

The MSE will balance these two effects, and the optimal value of  $\lambda$  (in terms of predictive performance) is one that achieves a good trade-off between bias and variance. This is the essence of the bias-variance trade-off in the context of ridge regression.

- Show that  $\beta_{\text{ridge}}^* = \frac{\beta_{\text{OLS}}^*}{1+\lambda}$  when  $X^T X = I_d$

**Answer:** The OLS estimator  $\beta_{\text{OLS}}^*$  is given by:

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

The ridge regression estimator  $\beta_{\text{ridge}}^*$  is given by:

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

Since  $X^T X = I_d$ , the OLS estimator simplifies to:

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

Now, considering the ridge regression estimator:

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

Since  $I_d + \lambda I_d$  is a diagonal matrix with each diagonal entry equal to  $1 + \lambda$ , its inverse is a diagonal matrix with each diagonal entry equal to  $\frac{1}{1+\lambda}$ . Thus, we have:

$$\beta_{\text{ridge}}^* = \frac{1}{1+\lambda} I_d X^T y$$

Since  $I_d X^T y$  is just  $X^T y$ , we obtain:

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

It looks like you've provided a description of the Elastic Net regularization method and its advantages over using Lasso or Ridge regularization individually. Here's the transcription of the content and the benefits of Elastic Net:

## Elastic Net

Using the previous notation, we can also combine Ridge and Lasso in the so-called Elastic Net regularization:

$$\beta_{\text{ENet}}^* = \arg \min_{\beta} \{ (y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1 \}$$

Calling  $\alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}$ , solving the previous Eq. is equivalent to:

$$\beta_{\text{ENet}}^* = \arg \min_{\beta} \{ (y_c - x_c \beta)^T (y_c - x_c \beta) + \lambda \left( \alpha \sum_{j=1}^d \beta_j^2 + (1 - \alpha) \sum_{j=1}^d |\beta_j| \right) \}$$

- This regularization overcomes some of the limitations of the Lasso, notably:
  - If  $d > N$  Lasso can select at most  $N$  variables  $\rightarrow$  ENet removes this limitation.
  - If a group of variables are highly correlated, Lasso randomly selects only one variable  $\rightarrow$  with ENet correlated variables have a similar value (grouped).
  - Lasso solution paths tend to vary quite drastically  $\rightarrow$  ENet regularizes the paths.
  - If  $N > d$  and there is high correlation between the variables, Ridge tends to have a better performance in prediction  $\rightarrow$  ENet combines Ridge and Lasso to have better (or similar) prediction accuracy with less (or more grouped) variables.

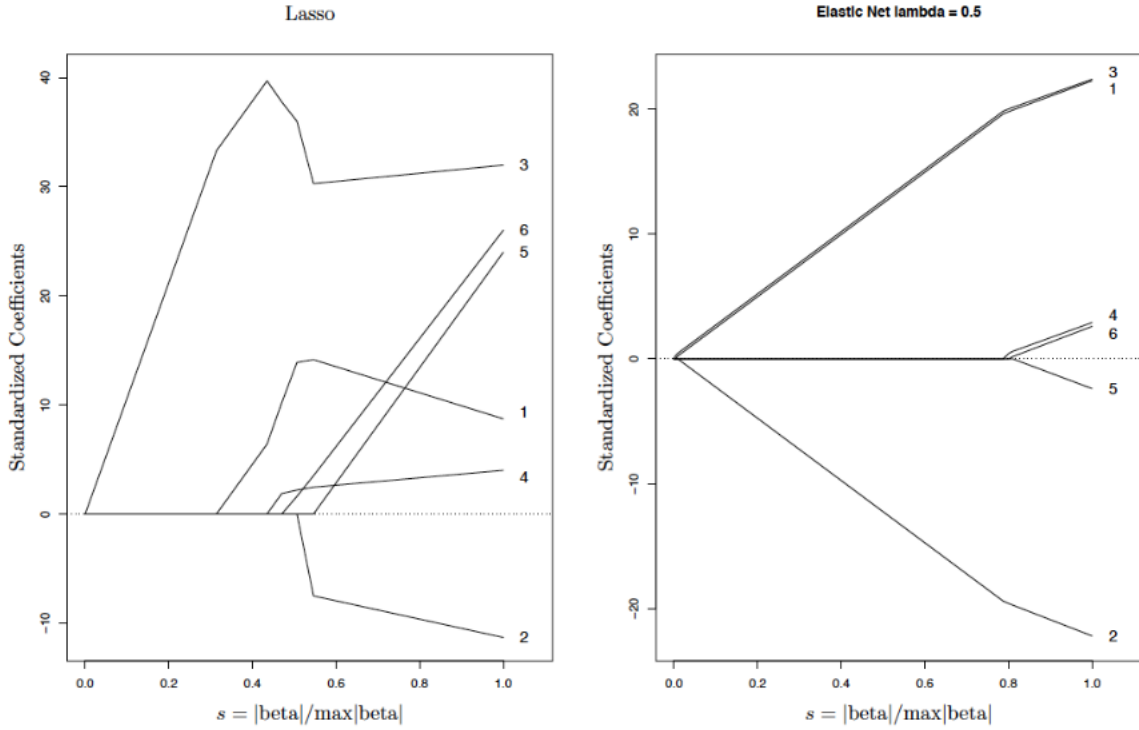


Figure 1: alt text

- Compute by hand the solution of Eq.2 supposing that  $X_c^T X_c = I_d$  and show that the solution is:

$$\beta_{\text{ENet}}^* = \frac{(\beta_{\text{OLS}}^*)_j \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$$

**Answer :**

To arrive at the Elastic Net solution using a thresholding approach, we start with the objective function given in Equation 2, taking into consideration that  $X_c^T X_c = I_d$  (the identity matrix):

$$\beta^{ENet} = \arg \min_{\beta} \{ (y_c - X_c \beta)^T (y_c - X_c \beta) + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \}$$

Because  $X_c^T X_c = I_d$ , the objective function simplifies to:

$$\beta^{ENet} = \arg \min_{\beta} \{ \|y_c - X_c \beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1 \}$$

The solution for the Ridge regression part (where  $\lambda_1 = 0$ ) with orthogonal predictors is:

$$\beta^{Ridge} = \frac{\beta^{OLS}}{1 + \lambda_2}$$

For Lasso regression, which uses an L1 penalty, we apply soft-thresholding to each coefficient. The soft-thresholding function for a given  $j$ -th coefficient, when  $X_c^T X_c = I_d$ , is defined as:

$$S_{\lambda_1}((\beta^{OLS})_j) = \text{sign}((\beta^{OLS})_j) (|(\beta^{OLS})_j| - \frac{\lambda_1}{2})_+$$

Here,  $(x)_+$  means  $\max(0, x)$ , and  $\text{sign}(x)$  is the sign function, which is  $+1$  for  $x > 0$ ,  $0$  for  $x = 0$ , and  $-1$  for  $x < 0$ .

In the Elastic Net, which combines both L1 and L2 penalties, the solution for each coefficient incorporates both the soft-thresholding from Lasso and the shrinkage from Ridge. The soft-thresholding operator is applied first, followed by the shrinkage due to the Ridge penalty:

$$\beta_j^{ENet} = \frac{S_{\lambda_1}((\beta^{OLS})_j)}{1 + \lambda_2}$$

Substituting the soft-thresholding function we get:

$$\beta_j^{ENet} = \frac{\text{sign}((\beta^{OLS})_j) (|(\beta^{OLS})_j| - \frac{\lambda_1}{2})_+}{1 + \lambda_2}$$

Now, we must account for the positive and negative scenarios depending on the sign of  $(\beta^{OLS})_j$ . If  $(\beta^{OLS})_j > \frac{\lambda_1}{2}$ , then  $\text{sign}((\beta^{OLS})_j) = +1$ , and if  $(\beta^{OLS})_j < -\frac{\lambda_1}{2}$ , then  $\text{sign}((\beta^{OLS})_j) = -1$ . If  $|(\beta^{OLS})_j| \leq \frac{\lambda_1}{2}$ , then the soft-thresholding output will be zero.

Thus, the final formula for each non-zero  $\beta_j^{ENet}$  is:

$$\beta_j^{ENet} = \frac{(\beta_j^{OLS}) \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$$

The  $\pm$  depends on the sign of the original OLS coefficient  $(\beta^{OLS})_j$ , which reflects the Lasso's characteristic of either subtracting or adding  $\frac{\lambda_1}{2}$  after thresholding, and then applying the Ridge shrinkage of  $\frac{1}{1+\lambda_2}$ .