

# Regularization

Machine Learning 2023-2024 - UMONS

Souhaib Ben Taieb

## Exercise 1

Consider the problem of multiple linear regression, where the aim is to find  $\hat{\beta}^{LS} = (\beta_0, \dots, \beta_p)^\top \in \mathbb{R}^{p+1}$  such that:

$$\hat{\beta}^{LS} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2.$$

Assuming that  $(\mathbf{X}^\top \mathbf{X})^{-1}$  is invertible, we can show that the ordinary least squares estimate is given by  $\hat{\beta}^{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  where  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{y} \in \mathbb{R}^n$ . Consider the problem of ridge regression, where the optimization problem is now formulated as finding  $\hat{\beta}^R = (\beta_0, \dots, \beta_p)^\top \in \mathbb{R}^{p+1}$  such that:

$$\hat{\beta}^R = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2,$$

where  $\lambda \geq 0$ . The solution is given by  $\hat{\beta}^R = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$ , where  $\mathbf{I}_p$  is the  $p \times p$  identity matrix. Let us consider a simple scenario where  $p = n$  and that  $\mathbf{X} = \mathbf{I}_p$ .

- Prove that  $\hat{\beta}^R = \frac{\hat{\beta}^{LS}}{\lambda + 1}$ .
- Given that  $\operatorname{Bias}(\hat{\beta}^{LS}) = 0$ , compute the bias of the ridge estimator.
- Assuming that the data generative process is  $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \in \mathbb{R}^n$  is a random noise vector, derive the covariance matrices of  $\hat{\beta}^{LS}$  and  $\hat{\beta}^R$  and show how they relate to one another. You can assume that  $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] = \sigma^2 \mathbf{I}_n$  and  $\mathbb{E}[\boldsymbol{\epsilon}] = 0$ .

The covariance matrix of a random vector  $\mathbf{a} \in \mathbb{R}^p$  is given by  $\operatorname{Cov}(\mathbf{a}) = \mathbb{E}[(\mathbf{a} - \mathbb{E}[\mathbf{a}])(\mathbf{a} - \mathbb{E}[\mathbf{a}])^\top] \in \mathbb{R}^{p \times p}$ .

**Solution :**

$$\begin{aligned} \hat{\beta}^{LS} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I}_p)^{-1} \mathbf{I}_p \mathbf{y} \\ &= \mathbf{y} \end{aligned}$$

$$\begin{aligned} \hat{\beta}^R &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{I}_p + \lambda \mathbf{I}_p)^{-1} \mathbf{I}_p \mathbf{y} \\ &= ((\lambda + 1) \mathbf{I}_p)^{-1} \mathbf{I}_p \mathbf{y} \\ &= \frac{1}{\lambda + 1} \mathbf{I}_p \mathbf{I}_p \mathbf{y} \\ &= \frac{\mathbf{y}}{\lambda + 1} \\ &= \frac{\hat{\beta}^{LS}}{\lambda + 1} \end{aligned}$$

The bias of the ridge estimator is obtained as:

$$\begin{aligned}
\text{Bias}(\hat{\beta}^R) &= \mathbb{E}[\hat{\beta}^R] - \beta \\
&= \mathbb{E}\left[\frac{\hat{\beta}^{LS}}{\lambda + 1}\right] - \beta \\
&= \frac{1}{\lambda + 1} \mathbb{E}[\hat{\beta}^{LS}] - \beta \\
&= \frac{\beta}{\lambda + 1} - \beta
\end{aligned}$$

Thus  $\hat{\beta}^R$  is a biased estimator of the true coefficients  $\beta$  if  $\lambda \neq 0$ .

The covariance matrix of the ordinary least squares coefficients is obtained as:

$$\begin{aligned}
\text{Cov}(\hat{\beta}^{LS}) &= \mathbb{E}\left[(\hat{\beta}^{LS} - \mathbb{E}[\hat{\beta}^{LS}])(\hat{\beta}^{LS} - \mathbb{E}[\hat{\beta}^{LS}])^\top\right] \\
&= \mathbb{E}\left[(\hat{\beta}^{LS} - \beta)(\hat{\beta}^{LS} - \beta)^\top\right] \\
&= \mathbb{E}\left[\hat{\beta}^{LS}(\hat{\beta}^{LS})^\top - \hat{\beta}^{LS}\beta^\top - \beta(\hat{\beta}^{LS})^\top + \beta\beta^\top\right] \\
&= \mathbb{E}\left[\hat{\beta}^{LS}(\hat{\beta}^{LS})^\top\right] - \mathbb{E}\left[\hat{\beta}^{LS}\beta^\top\right] - \mathbb{E}\left[\beta(\hat{\beta}^{LS})^\top\right] + \mathbb{E}\left[\beta\beta^\top\right] \\
&= \mathbb{E}[\mathbf{y}\mathbf{y}^\top] - \beta\beta^\top - \beta\beta^\top + \beta\beta^\top \\
&= \mathbb{E}\left[(\mathbf{X}\beta + \boldsymbol{\epsilon})(\mathbf{X}\beta + \boldsymbol{\epsilon})^\top\right] - \beta\beta^\top \\
&= \mathbb{E}[\mathbf{X}\beta\beta^\top\mathbf{X}^\top] + \mathbb{E}[\mathbf{X}\beta\boldsymbol{\epsilon}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\beta^\top\mathbf{X}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] - \beta\beta^\top \\
&= \beta\beta^\top + \sigma^2\mathbf{I}_p - \beta\beta^\top \\
&= \sigma^2\mathbf{I}_p
\end{aligned}$$

The covariance matrix of the ridge coefficients are given by:

$$\begin{aligned}
\text{Cov}(\hat{\beta}^R) &= \mathbb{E}\left[(\hat{\beta}^R - \mathbb{E}[\hat{\beta}^R])(\hat{\beta}^R - \mathbb{E}[\hat{\beta}^R])^\top\right] \\
&= \mathbb{E}\left[\left(\frac{\hat{\beta}^{LS}}{\lambda + 1} - \frac{\beta}{\lambda + 1}\right)\left(\frac{\hat{\beta}^{LS}}{\lambda + 1} - \frac{\beta}{\lambda + 1}\right)^\top\right] \\
&= \frac{1}{(\lambda + 1)^2} \left( \mathbb{E}[\hat{\beta}^{LS}(\hat{\beta}^{LS})^\top] - \mathbb{E}[\hat{\beta}^{LS}\beta^\top] - \mathbb{E}[\beta(\hat{\beta}^{LS})^\top] + \mathbb{E}[\beta\beta^\top] \right) \\
&= \frac{1}{(\lambda + 1)^2} (\sigma^2\mathbf{I}_p + \beta\beta^\top - \beta\beta^\top - \beta\beta^\top + \beta\beta^\top) \\
&= \frac{\sigma^2\mathbf{I}_p}{(\lambda + 1)^2} \\
&= \frac{\text{Cov}(\hat{\beta}^{LS})}{(\lambda + 1)^2}
\end{aligned}$$

## Exercise 2

Suppose that the columns of  $\mathbf{X}_1$  are orthonormal, and that  $\mathbf{X}_2 = 10\mathbf{X}_1$ . Show that the ordinary least squares estimates are equivariant, meaning that multiplying  $\mathbf{X}$  by a constant  $c$  scales the coefficients estimates by a factor  $\frac{1}{c}$ . Is it also the case for the ridge estimates?

**Solution :**

The least square coefficients from  $\mathbf{X}_1$  would be  $\hat{\beta}_1^{LS} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}$ . From  $\mathbf{X}_2$ , we would have:

$$\begin{aligned}\hat{\beta}_2^{LS} &= (\mathbf{X}_2^T \mathbf{X}_2)^{-1} \mathbf{X}_2^T \mathbf{y} \\ &= (100\mathbf{X}_1^T \mathbf{X}_1)^{-1} 10\mathbf{X}_1^T \mathbf{y} \\ &= \frac{1}{10} (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y} \\ &= \frac{1}{10} \hat{\beta}_1^{LS}.\end{aligned}$$

Similarly, the ridge coefficients obtained from  $\mathbf{X}_1$  would be  $\hat{\beta}_1^R = \frac{1}{1+\lambda} \mathbf{X}_1^T \mathbf{y}$ . From  $\mathbf{X}_2$ , we would have:

$$\begin{aligned}\hat{\beta}_2^R &= (\mathbf{X}_2^T \mathbf{X}_2 + \lambda \mathbf{I}_p)^{-1} \mathbf{X}_2^T \mathbf{y} \\ &= (100\mathbf{X}_1^T \mathbf{X}_1 + \lambda \mathbf{I}_p)^{-1} 10\mathbf{X}_1^T \mathbf{y} \\ &= (100\mathbf{I}_p + \lambda \mathbf{I}_p)^{-1} 10\mathbf{X}_1^T \mathbf{y} \\ &= \frac{10}{100 + \lambda} \mathbf{I}_p^{-1} \mathbf{X}_1^T \mathbf{y} \\ &= \frac{10}{100 + \lambda} \mathbf{X}_1^T \mathbf{y} \\ &\neq \frac{1}{10} \hat{\beta}_1^R.\end{aligned}$$

The ridge estimates are not equivariant, meaning that they are sensitive to the scale of the data.