

FYS-STK4155 Week 37

Janita Ovidie Sandtrøen Willumsen

(Dated: September 17, 2023)

EXERCISE 1

We have assumed that our data can be described by the continuous function $f(\mathbf{x})$, and an error term $\epsilon \sim N(0, \sigma^2)$. If we approximate the function with the solution derived from a model $\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta}$ the data can be described with $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$. The expectation value

$$\begin{aligned}\mathbb{E}(\mathbf{y}) &= \mathbb{E}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \mathbb{E}(\mathbf{X}\boldsymbol{\beta}) + \mathbb{E}(\boldsymbol{\epsilon}) && \text{where the expected value } \boldsymbol{\epsilon} = \mathbf{0} \\ \mathbb{E}(y_i) &= \sum_{j=0}^{P-1} X_{i,j} \beta_j && \text{for the each element} \\ &= X_{i,*} \beta_i && \text{where } * \text{ replace the sum over index } i\end{aligned}$$

The variance for the element y_i can be found by

$$\begin{aligned}\mathbb{V}(y_i) &= \mathbb{E}[(y_i - \mathbb{E}(y_i))^2] \\ &= \mathbb{E}(y_i^2) - (\mathbb{E}(y_i))^2 \\ &= \mathbb{E}((X_{i,*} \beta_i + \epsilon_i)^2) - (X_{i,*} \beta_i)^2 \\ &= \mathbb{E}((X_{i,*} \beta_i)^2 + 2\epsilon_i X_{i,*} \beta_i + \epsilon_i^2) - (X_{i,*} \beta_i)^2 \\ &= \mathbb{E}((X_{i,*} \beta_i)^2) + \mathbb{E}(2\epsilon_i X_{i,*} \beta_i) + \mathbb{E}(\epsilon_i^2) - (X_{i,*} \beta_i)^2 \\ &= (X_{i,*} \beta_i)^2 + \mathbb{E}(\epsilon_i^2) - (X_{i,*} \beta_i)^2 \\ &= \mathbb{E}(\epsilon_i^2) = \sigma^2\end{aligned}$$

The expression for the optimal parameter

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

We find the expected value of $\hat{\boldsymbol{\beta}}$

$$\begin{aligned}\mathbb{E}(\hat{\boldsymbol{\beta}}) &= \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y}) && \text{using that } \mathbf{X} \text{ is a non-stochastic variable} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} && \text{using } \mathbb{E}(\mathbf{y}) = \mathbf{X} \boldsymbol{\beta} \\ &= \boldsymbol{\beta}\end{aligned}$$

we can find the variance by

$$\begin{aligned}
\mathbb{V}(\hat{\beta}) &= \mathbb{E}[(\hat{\beta} - \mathbb{E}(\hat{\beta}))^2] \\
&= \mathbb{E}(\hat{\beta}\hat{\beta}^T) - \mathbb{E}(\hat{\beta})\mathbb{E}(\hat{\beta})^T \\
&= \mathbb{E}(((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})^T) - \hat{\beta}\hat{\beta}^T \\
&= \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) - \hat{\beta}\hat{\beta}^T \\
&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y} \mathbf{y}^T) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \hat{\beta}\hat{\beta}^T \\
&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta \beta^T \mathbf{X}^T + \sigma^2 \mathbf{I}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \hat{\beta}\hat{\beta}^T \\
&= \beta \beta^T + \sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) - \hat{\beta}\hat{\beta}^T \\
&= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}
\end{aligned}$$

Knowing the expectation value and the variance of $\hat{\beta}$, we can define a confidence interval for each $\hat{\beta}_j \pm \text{std}(\hat{\beta}_j)$ for $j = 1, 2, \dots, P-1$.

EXERCISE 2

Last week we showed that the optimal $\hat{\beta}^{Ridge}$ can be derived from MSE, and is defined as

$$\hat{\beta}^{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

The expectation value is then

$$\begin{aligned}
\mathbb{E}(\hat{\beta}^{Ridge}) &= \mathbb{E}((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}) \\
&= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y}) && \text{since } \mathbf{X} \text{ and } \lambda \mathbf{I} \text{ are non-stochastic variables} \\
&= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} \beta && \text{using } \mathbb{E}(\mathbf{y}) \text{ from exercise 1}
\end{aligned}$$

For $\lambda = 0$ we have $\mathbb{E}(\hat{\beta}^{OLS})$. The variance

$$\begin{aligned}
\mathbb{V}(\hat{\beta}^{Ridge}) &= \mathbb{E}(\hat{\beta}_R \hat{\beta}_R^T) - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= \mathbb{E}(((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y})((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y})^T) - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= \mathbb{E}((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \mathbf{y}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1})^T) - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y} \mathbf{y}^T) \mathbf{X} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1})^T - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{X} \beta \beta^T \mathbf{X}^T + \sigma^2 \mathbf{I}) \mathbf{X} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1})^T - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1})^T + (\mathbb{E}(\hat{\beta}_R))^2 - (\mathbb{E}(\hat{\beta}_R))^2 \\
&= \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1})^T
\end{aligned}$$