

Lecture 16: Ridge Regression

• Linear regression based on a version of $\hat{\mu}^{MLE}$. we usually assume an n -dim vector $y = (y_1, \dots, y_n)^T$ from linear model

$$y = X \cdot \beta + \epsilon$$

$X^{n \times p}$ known as design/structure matrix; β is an unknown p -dimensional parameter vector & noise vector ϵ is assumed to have uncorrelated components. This $\epsilon \sim (0, \sigma^2 \cdot I_n)$ & constant variance

I identity matrix of size n



The least squares estimator (Gauss & Legendre in the early 1800s!) minimizes SSE

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \|y - X\beta\|^2 \}$$

$$-2(y - X\beta) \cdot X = 0$$

$p \times 1$ $n \times 1$

$$(y - X\beta) \cdot X = 0$$

$$X^T y - X^T X \cdot \beta = 0$$

$$X^T X \beta = X^T y \quad ; \quad \hat{\beta} = (X^T X)^{-1} \cdot X^T y$$

$$\text{Let } S = (X^T X) \Rightarrow \hat{\beta} = S^{-1} X^T y$$

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \{ \|y - X \cdot \beta\|^2 + \lambda \cdot \| \beta \|^2 \}$$

$$\text{Recall } \| \beta \|^2 = \beta_1^2 + \beta_2^2 + \dots + \beta_p^2$$

Then, do similar math

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

$$-X^T y + X^T X \cdot \beta + \lambda \beta = 0$$

$$\beta (X^T X + \lambda \cdot I_n) = X^T y$$

$$\hat{\beta}^{\text{Ridge}} = (X^T X + \lambda \cdot I_n)^{-1} \cdot X^T y$$

Standard errors on $\hat{\beta}$

$$\hat{\beta}^{OLS} \sim (\beta, \sigma^2 \cdot S^{-1}) \quad \Bigg| \quad \hat{\beta}^{ridge} \sim ((S + \lambda \cdot I_n)^{-1} \cdot S \cdot \beta,$$



$$\sigma^2 \cdot (S + \lambda \cdot I_n)^{-1} \cdot S \cdot (S + \lambda \cdot I_n)^{-1}$$