

Lecture Erasmus+
course, October 16,
2023

Ordinary least squares (OLS)

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \sum_j x_{ij} \beta_j)^2$$

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} C(\beta) (= MSE)$$

Design matrix $\hat{X} = \begin{bmatrix} 1 & x_{00} & \dots & x_{0p} \\ \vdots & \vdots & & \vdots \\ x_{n-1,0} & \dots & x_{n-1,p} \end{bmatrix}$

$$\hat{X} \in \mathbb{R}^{n \times p}$$

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

in case $\hat{X}^T \hat{X}$ is singular

$$\hat{X}^T \hat{X} + \lambda I_{p \times p} \quad \lambda \geq 0$$

$$\hat{X}^T \hat{X} \in \mathbb{R}^{p \times p}$$

with MSE as $C(\beta)$

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = \frac{2}{m} X^T X$$

Hessian
matrix

using SVD

$$U U^T = U^T U = \mathbb{I} = V^T V = V V^T$$

$$U \in \mathbb{R}^{n \times n}$$

$$V \in \mathbb{R}^{P \times P}$$

$$\Sigma = \begin{bmatrix} \sigma_0 & & & \\ & \sigma_1 & & \\ & & \ddots & \\ & & & \sigma_{P-1} & 0 \end{bmatrix}$$

$$\sigma_0 > \sigma_1 > \sigma_2 > \dots > \sigma_{P-1} > 0$$

$$X^T X = V \Sigma \underbrace{U^T U}_{P \times P} \Sigma V^T$$

$$= V \underbrace{\Sigma \Sigma}_{P \times P} V^T = V \Sigma^2 V^T$$

$$\Sigma^2 = \begin{bmatrix} \sigma_0^2 & & \\ & \ddots & \\ & & \sigma_{P-1}^2 \end{bmatrix}$$

$$(X^T X) V = V \Sigma^2$$

$$V = \begin{bmatrix} v_0 & v_1 & \dots & v_{P-1} \end{bmatrix}$$

$$(X^T X) v_i = \sigma_i^2 v_i^T v_j = S_{ij}$$

$$X^T X + \lambda \mathbb{I} \Rightarrow \lambda \geq 0$$

Ridge Regression

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \sum_j x_{ij} \beta_j)^2$$

$$+ \lambda \sum_{j=0}^{p-1} \beta_j^2$$

it is common to leave out
 β_0 ($j=0$), the intercept.

$$\frac{\partial C}{\partial \beta} = 0 \Rightarrow \hat{\beta}_{\text{Ridge}} = (X^T X + \lambda \mathbb{I})^{-1} X^T y$$

$$MSE = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2$$

$$\hat{y} = X \hat{\beta}$$

$$C(\beta)_{Ridge} = \frac{1}{n} \| (y - X\beta) \|_2^2$$

$$+ \lambda \| \beta \|_2^2$$

(common to leave out β_0)

$$C(\beta)_{Lasso} = \frac{1}{n} \| (y - X\beta) \|_2^2$$

$$+ \lambda \| \beta \|_1$$

$$\| \beta \|_1 = \sum_{j=0}^{p-1} |\beta_j|$$

Lasso regression we have

$$\frac{d |\beta_j|}{d \beta_j} = \begin{cases} +1 & \beta_j > 0 \\ -1 & \beta_j < 0 \end{cases}$$

\Rightarrow optimization with constraint
(no analytical solution)

$$\hat{\beta}_{OLS} = \left[\sum_{j=0}^{p-1} u_j u_j^T \right] y$$

$$\hat{\beta}_{Ridge} = \left[\sum_{j=0}^{p-1} \frac{u_j^2}{\sigma_j^2 + \lambda} u_j u_j^T \right] y$$

with a given λ , we can shrink the role of a specific parameter β_j

$$X = \begin{bmatrix} 2 & 0 \\ c_1 & 1 \\ c_2 & c \end{bmatrix}$$

$$\overset{\lambda}{\beta}_{\text{par}} = \begin{bmatrix} ? \\ ? \end{bmatrix}$$

$$y = \begin{bmatrix} 2 & 0 \\ c_1 & 1 \\ c_2 & c \end{bmatrix} \begin{bmatrix} ? \\ ? \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \\ 0 \end{bmatrix}$$

$$(y = \begin{bmatrix} 4 \\ 2 \\ 3 \end{bmatrix})$$

$$MSE = \frac{1}{n=3} \sum_{i=0}^2 (y_i - \hat{y}_i)^2 = 3$$

Basic assumptions

$$y = f(x) + \varepsilon$$

\uparrow Random noise

non-stochastic and
continuous function

$$\varepsilon \sim N(0, \sigma^2)$$

$$p(y) = ?$$

$$E[y] = ?$$

$$E[y] = \int_{y \in D} p(y) y dy$$

$$\left\{ \begin{aligned} &= \sum_{y_i \in D} p(y_i) y_i \end{aligned} \right.$$

in most cases we do not know
 $P(g)$.

sample mean

$$E[g] = \frac{1}{n} \sum_{i=0}^{n-1} g_i \neq \text{True}$$

$E[S]$

$$= \int_{g \in D} P(g) g dg$$

$$E[S] = E[f(x)] + E[\varepsilon] = 0$$

$$E[f(x)] = \int P(g) f(x) dx =$$

$$f(x) \underbrace{\int P(g) dx}_{=1}$$

$$= f(x)$$

$$E[y] = f(x) \simeq x\beta$$

$$E[y_i] = \sum_{j=0}^{p-1} x_{ij} \beta_j = x_i * \beta$$

$$\begin{aligned} \text{var}[y_i] &= E[(y_i - E[\bar{y}_i])^2] \\ &= E[\bar{y}_i^2] - (E[\bar{y}_i])^2 \\ &= E[(x_i * \beta + \varepsilon_i)^2] - (x_i * \beta)^2 \\ &= E[(x_i * \beta)^2 + 2x_i * \beta \varepsilon_i + \varepsilon_i^2] \\ &\quad - (x_i * \beta)^2 \quad \stackrel{2x_i * \beta E(\varepsilon_i)}{=} 0 \end{aligned}$$

$$= E[\epsilon_i^2] = \text{var}[\epsilon_i] = \sigma^2$$

we can assume that

$$y_i \sim N(x_{i*}\beta, \sigma^2)$$

mean value

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{(y_i - x_{i*}\beta)^2}{2\sigma^2} \right]$$

(From Bayes' theorem)

\Rightarrow Derivation of OLS

$$y_i \in D = \{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_n)\}$$

y_i are independent and identically distributed (i.i.d.)

$$P(\underbrace{x, y}_D | \beta) = \prod_{i=0}^{n-1} P(y_i | x_i, \beta)$$

$$= \prod_{i=0}^{n-1} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - x_i * \beta)^2}{2\sigma^2}\right]$$

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^p} P(D | \beta)$$

MLE

Max-likelihood estimator

equivalent problem

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^P} \log P(D|\beta)$$

or

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} -\log P(D|\beta)$$

$$C(\beta) = -\log P(D|\beta)$$

$$= - \sum_{i=0}^{n-1} \log P(y_i | x_i | \beta)$$

$$= \frac{n}{2} \log (2\pi\sigma^2) +$$

$$+ \sum_{i=0}^{n-1} \frac{(y_i - x_i * \beta)^2}{2\sigma^2}$$

$$= \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|(\mathbf{y} - \mathbf{x}\beta)\|_2^2$$

$$\frac{\partial C}{\partial \beta} = 0 = -\mathbf{x}^T(\mathbf{x}\beta - \mathbf{y}) = 0$$

$$\Rightarrow \boxed{\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}}$$

$$y = f(x) + \varepsilon$$

For Ridge & Lasso, we need
Bayer's theorem

$$P(\beta | D) \propto P(D | \beta) P(\beta)$$

posterior distribution

likelihood

prior