

Lecture January 19

Dataset $\{(x_0, y_0), (x_1, y_1), \dots, (x_{n-1}, y_{n-1})\}$

Assumption

$$y(x) = f(x) + \underset{\substack{\uparrow \\ \text{normally} \\ \text{distributed}}}{\varepsilon}$$

$$f(x) \approx \tilde{y} = \text{Model } \varepsilon \sim N(0, \sigma^2)$$

$$\tilde{y} \in \mathbb{R}^n \quad y \in \mathbb{R}^n$$

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

$$\beta^T = [\beta_0 \ \beta_1 \ \dots \ \beta_{p-1}]$$

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

Error function (cost/loss--)

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

$$= \frac{1}{n} \sum_i (y_i - X_i \beta)^2$$

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

$$\frac{\partial C(\beta)}{\partial \beta_j} = 0$$

$$\beta^{\text{opt}} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} C(\beta)$$

$$X = \begin{bmatrix} x_{00} & x_{01} & \dots & x_{0p-1} \\ x_{10} & & & \\ \vdots & & & \\ x_{n-1,0} & \dots & \dots & x_{n-1,p-1} \end{bmatrix}$$

$$\frac{\partial C(\beta)}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} \left(\frac{1}{n} \sum_{i=0}^{n-1} \right)$$

$$\left(y_i - \beta_0 x_{i0} - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots \right)^2$$

$$= -\frac{2}{n} \sum_{i=0}^{n-1} \left(\underline{x_{ij}} [y_i - \beta_0 x_{i0} - \beta_1 x_{i1} - \dots] \right)$$

Matrix-vector form ($\frac{2}{n}$)

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

$$\Rightarrow X^T y = X^T X \beta^{opt} \Rightarrow$$

$$\boxed{\underline{\beta}^{opt} = (\underline{X^T X})^{-1} \underline{X^T y}} \quad \begin{array}{l} \text{RHS} \\ \text{is} \\ \text{known} \end{array}$$

ORDINARY LEAST SQUARES

$$\hat{y} = \hat{y}_{predict} = X \beta^{opt}$$

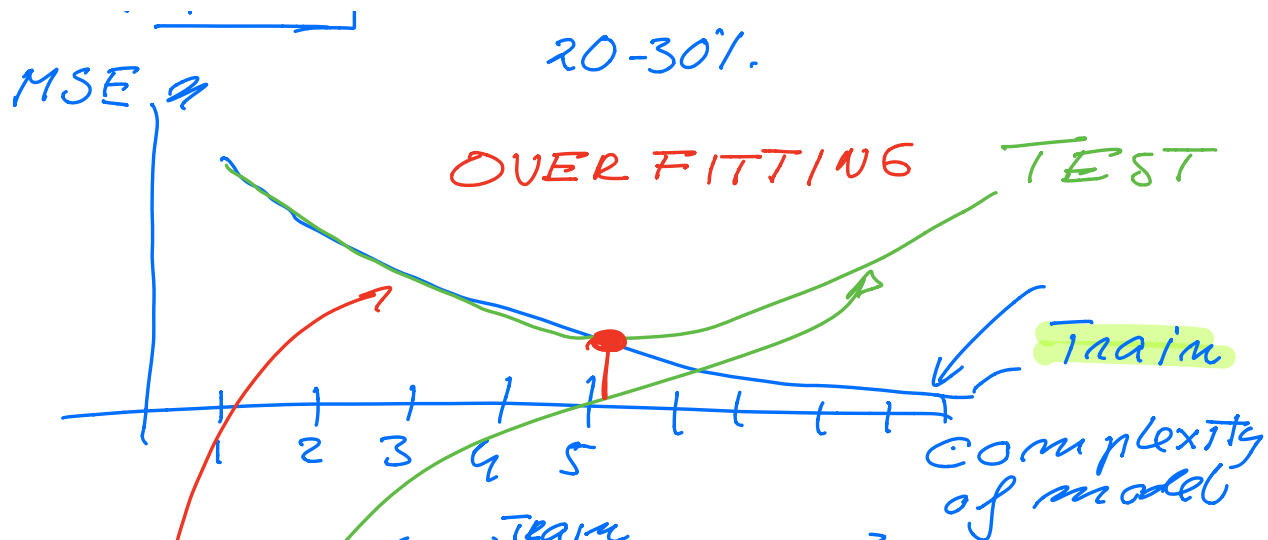
splitting data in train
(validation), test

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

TRAIN

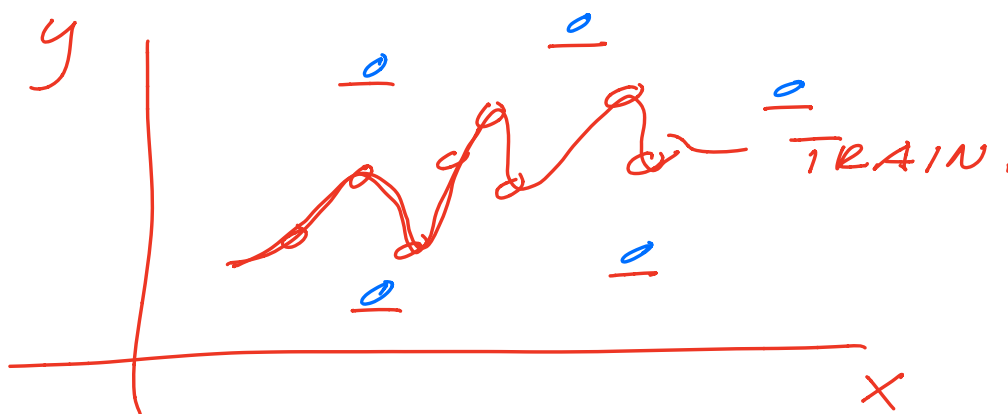
70-80% of all
data - n -

$$\hat{y}_{\text{test}} = X_{\text{test}}^T \beta_{OLS}^{opt}$$



$$MSE_{\text{Train}} = \frac{1}{n_{\text{Train}}} \sum_{i=1}^{n_{\text{Train}}} (y_i - \tilde{y}_i'(\text{Train}))^2 \quad (\text{poly-degree})$$

$$MSE_{\text{Test}} = \frac{1}{n_{\text{Test}}} \sum_{i=1}^{n_{\text{Test}}} (y_i - \tilde{y}_i'(\text{Test}))^2$$



Lasso & Ridge Regression's

OLS

$$C(\beta) = \frac{1}{n} [(y - X\beta)(y - X\beta)]$$

norm-2 of a vector;

$$\|x\|_2 = \sqrt{\sum_i x_i^2}$$

$$\|x\|_2^2 = \sum_i x_i^2$$

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

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

$$\hat{y} = X \beta^{\text{opt}}$$

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

$$\beta \in \mathbb{R}^p$$

$(X^T X)^{-1}$ can be singular.

Ridge

$A =$

$$\begin{bmatrix} a_{00} & a_{01} & \dots & a_{0n} \\ \vdots & a_{11} & & \\ \vdots & & a_{22} & \\ a_{n0} & & & \dots & a_{nn} \end{bmatrix}$$

$$\tilde{A} = A + \lambda I$$

$$= \begin{bmatrix} a_{00} + \lambda & a_{01} & \dots & a_{0n} \\ \vdots & a_{11} + \lambda & & \\ a_{n0} & & & a_{nn} + \lambda \end{bmatrix}$$

$$\lambda \sim 10^{-5}$$

$$(X^T X)^{-1} \rightarrow (X^T X + \lambda I)^{-1}$$

"Ridge" :

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

Ridge - cost

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

$$\lambda \geq 0$$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \Rightarrow$$

$$\boxed{\beta_{\text{Ridge}}^{\text{opt}} = (X^T X + \lambda I)^{-1} X^T y}$$

Lasso

$$\|x\|_1 = \sum_i |x_i|$$

$$C(\beta) = \frac{1}{n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

can reduce MSE compared
to OLS by tuning λ

hyperparameter