

Lecture Notes September 17

OLS, Ridge & Lasso

$$\text{OLS} : \|(y - x\beta)\|_2^2.$$

$$\Rightarrow \frac{1}{n} \sum_{i=0}^{n-1} (y_i - x_i \cdot \beta)(y_i - x_i \cdot \beta)$$

$$\text{Ridge} : \|(y - x\beta)\|_2^2 + \lambda \|\beta\|_2^2$$

$$\|\beta\|_2^2 = \sum_{i=0}^{p-1} \beta_i^2 \leq t$$

Lasso Regression:

$$\|(y - x\beta)\|_2^2 + \lambda \|\beta\|_1$$

$$\|\beta\|_1 = \sum_{i=0}^{p-1} |\beta_i| \leq t$$

$$\lambda \geq 0$$

$$\hat{\beta}^{\text{OLS}} = (x^T x)^{-1} x^T y$$

$$\hat{\beta}^{\text{Ridge}} = (x^T x + \lambda \mathbb{I})^{-1} x^T y$$

$\hat{\beta}^{\text{Lasso}}$: not analytical solution

$$\frac{d|x|}{dx} = \text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \end{cases}$$

$$\frac{\partial C^{\text{Lasso}}(\beta)}{\partial \beta} = -2x^T(y - x\beta) + \lambda \text{sign}(\beta)$$

(-1 if $x < 0$)

$$\frac{\partial C^{\text{Ridge}}}{\partial \beta} = -2x^T(y - x\beta) + 2\beta$$

$$\frac{\partial C^{\text{OLS}}}{\partial \beta} = -2x^T(y - x\beta)$$

$$\|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$$

$$\|v\|_2^2 = \sum v_i^2$$

$$\|v\|_1 = \sum_i |v_i|$$

$$\frac{\partial^2 C^{\text{OLS}}}{\partial \beta \partial \beta^T} = 2x^T x = 2 \cdot \textcircled{H}$$

Hessian

$$\text{var}(\hat{\beta}^{\text{OLS}}) = \sigma^2 (x^T x)^{-1}$$

$$\text{var}(\hat{\beta}_i^{\text{OLS}}) = \sigma^2 (x^T x)_{ii}^{-1}$$

$$\overset{\sim}{\epsilon} \sim N(0, \sigma^2)$$

Covariance matrix $= C[\bar{x}]$

$$= \frac{1}{n} \bar{X}^T \bar{X}$$

$$= \begin{bmatrix} \text{var}[x_0] & \text{cov}[\bar{x}_0 \bar{x}_1] & \dots & \text{cov}[\bar{x}_0 \bar{x}_{p-1}] \\ \vdots & \ddots & & \\ \text{cov}[\bar{x}_{p-1} \bar{x}_0] & \dots & \dots & \text{var}[\bar{x}_{p-1}] \end{bmatrix}$$

$$\text{var}[x_0] = \frac{1}{n} \sum_{i=0}^{n-1} (x_{i0} - \bar{x}_0)^2$$

$\hat{A} \propto C[\bar{x}]$ a positive definite matrix

$\Rightarrow \hat{A} > 0 \Rightarrow \hat{\beta}$ is a true minimum of $C(\beta)$

$$\frac{\partial^2 C_{\text{Ridge}}}{\partial \beta \partial \beta^T} = 2 \boxed{H} + 2\lambda \beta$$

$$\frac{\partial^2 C_{\text{Lasso}}}{\partial \beta \partial \beta^T} = 2H$$

- - -

$$\frac{\partial C}{\partial \beta}^{\text{lasso}} = -2X^T y + 2H\beta + \lambda \text{sign}(\beta)$$

$C[x]$ = covariance matrix

$$C(\beta) = C(y|x\beta) =$$

cost function

$$= L(y|x\beta)$$

loss function

$$= R(y|x\beta)$$

Risk function,

$$\begin{cases} \text{var}[x] & |E[x] \\ \text{cov}[x,y] & \text{statistical} \\ & \text{properties.} \end{cases}$$

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

$$X = U \Sigma V^T$$

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

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

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

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

$$UU^T = U^T U = I$$

$$V^T V = V V^T = I$$

$$(XX^T)U = U\Sigma\Sigma^T = U\Sigma^2$$

$$\hat{y}_{OLS} = X\hat{\beta}_{OLS} = \frac{\sum_{i=0}^{p-1} u_i u_i^T y}{u_i u_i^T}$$

$$\hat{y}_{Ridge} = X\hat{\beta}_{Ridge} = \frac{\sum_{i=0}^{p-1} u_i \frac{\tau_i^2}{\tau_i^2 + \lambda} u_i^T y}{\sum_{i=0}^{p-1} \frac{\tau_i^2}{\tau_i^2 + \lambda}}$$

$$\tau_i \geq \tau_{i+1}$$

Ridge shrinks the singular values, reducing more (for a given λ) those τ_i which are smaller. These correspond to eigenvalues with small variance.

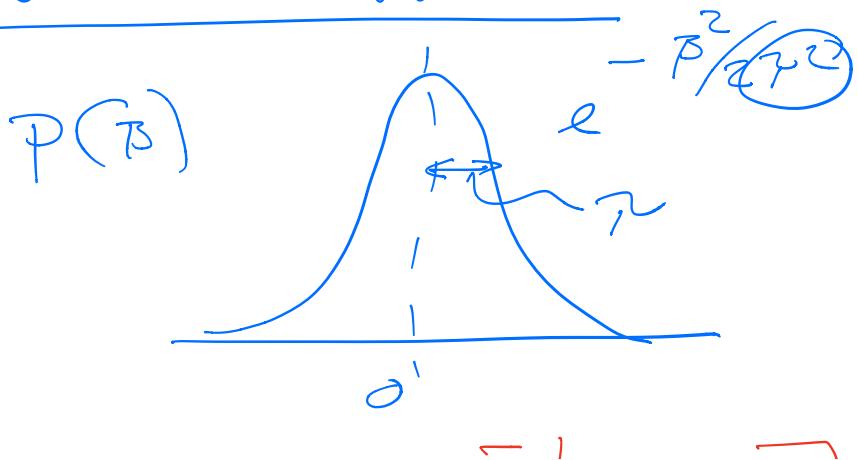
$$\begin{aligned}
 \hat{\beta}^{\text{OLS}} &= \frac{(X^T X)^{-1} X^T y}{\sigma^2 \Sigma^T u^T u \Sigma} \\
 &= (\Sigma^T u^T u \Sigma)^{-1} \Sigma^T u^T y \\
 &= (\Sigma^T \Sigma)^{-1} \Sigma^T u^T y
 \end{aligned}$$

|| $\text{var}[\hat{\beta}^{\text{Ridge}}] = \sigma^2 [X^T X + \lambda I]^{-1}$

$\times X^T X [X^T X + \lambda I]$

Ridge shrinks the variance of $\hat{\beta}_i^{\text{Ridge}}$ more for those β_i values with small singular values.

Bayesian approach



$$X^T X = \mathbb{I} = \begin{bmatrix} 1 & 0 \\ 0 & \ddots \\ 0 & 0 \end{bmatrix}$$

$$\hat{y}_{\text{Ridge}} = \frac{1}{1+\lambda} \hat{y}_{\text{OLS}}$$

$$\frac{X^T X = \mathbb{I}}{X^T X + \mathbb{I}} = \frac{1}{1+\lambda} y$$

$$\text{var}[\hat{\beta}_{\text{Ridge}}] = \sigma^2 \frac{1}{\mathbb{I}(1+\lambda)} \frac{\pi}{\mathbb{I}(1+\lambda)} \\ = \sigma^2 \frac{1}{(1+\lambda)^2}$$

Degrees of freedom

P - predictors/features - - .

$$\hat{y}_{\text{OLS}} = \sum_{i=0}^{P-1} u_i u_i^T y$$

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

$$\text{Tr}[X(X^T X)^{-1} X^T] = P$$

= degrees of freedom,

$$\hat{y}_{\text{Ridge}} = X(X^T X^{-1} \lambda I)^{-1} X^T y$$

$$\text{Tr} [x(x^T x + \lambda \mathbb{I})^{-1} x^T] = \sum_{j=0}^{P-1} \frac{\sigma_j^2}{\sigma_j^2 + \lambda} = df(\lambda)$$

$\lambda = 0 \Rightarrow df(\lambda) = P$

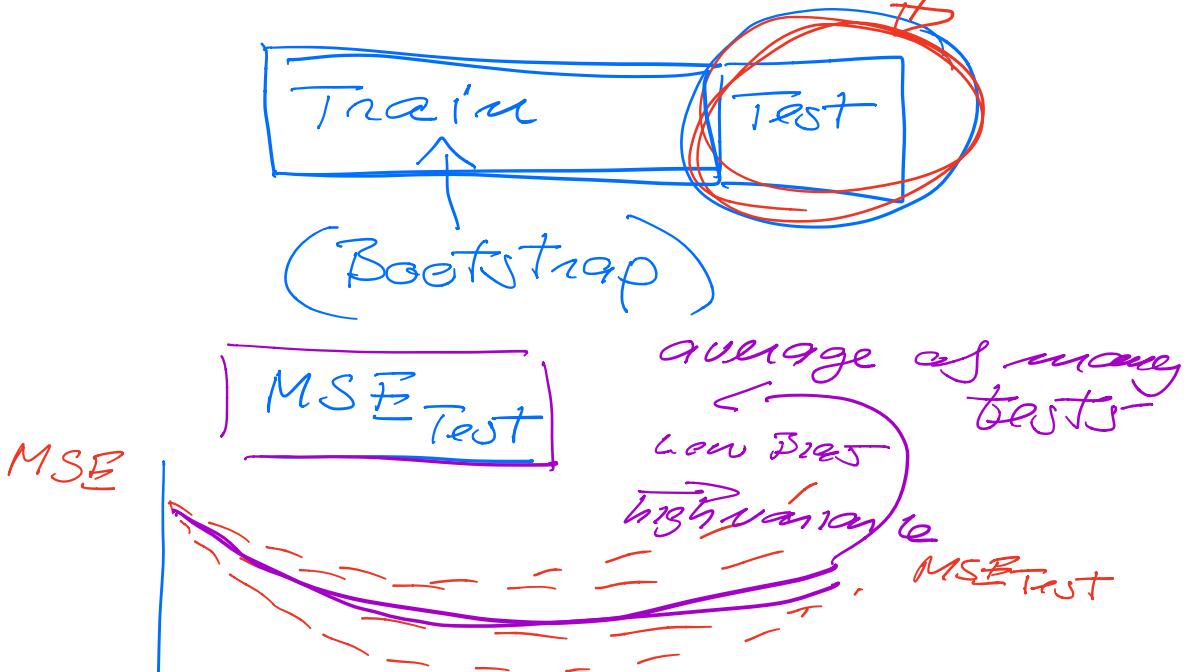
$$\text{Tr}[A] = \sum_{i=1}^n a_{ii}$$

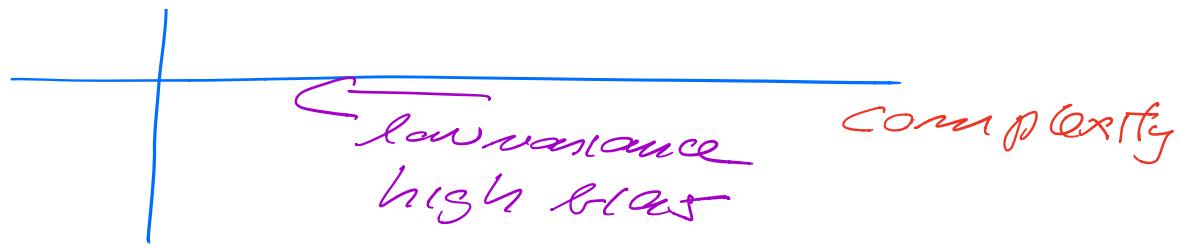
$$\lambda \rightarrow \infty \Rightarrow df(\lambda) = 0$$

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

Resampling methods —

1b : Bootstrap + bias-variance
+ 1d (OLS + Ridge)

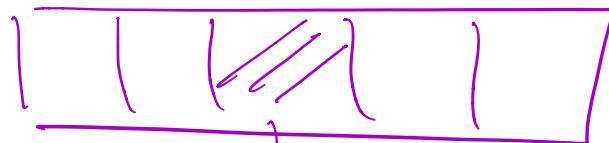
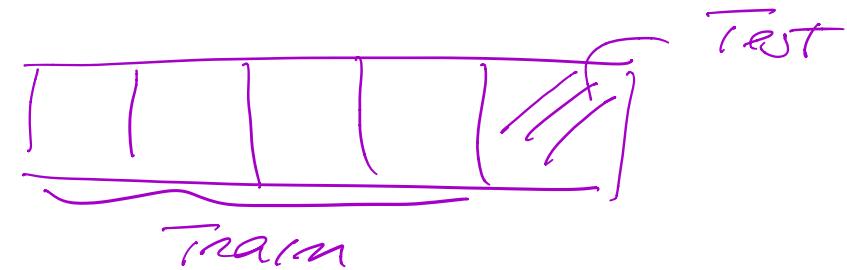




IC

OLS + Cross-validation

K-Folds $K = 5$



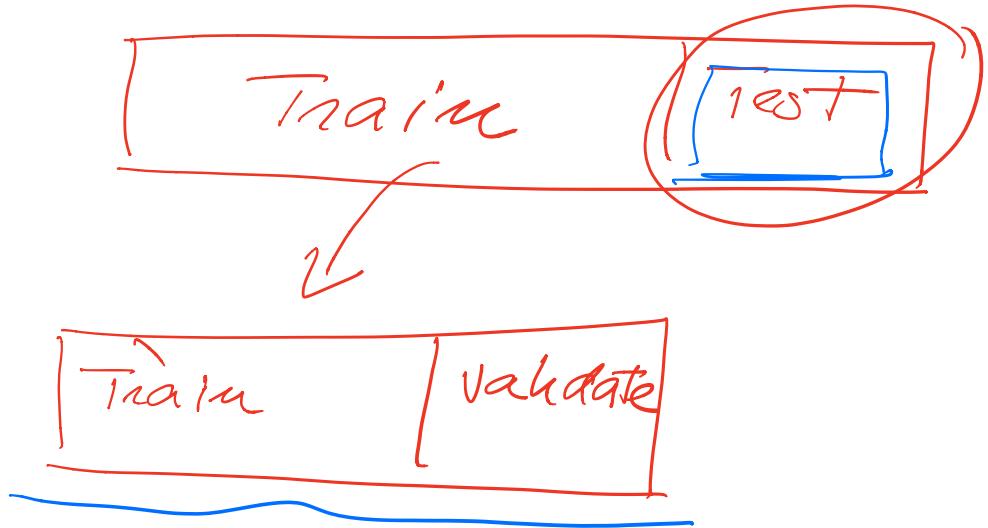
↓

3 more times

in total 5 MSE_{test}
calculations.

- No Need to do a Train-Test split.

Train-test split



Cross-validate here.

Ridge : polynomial complexity
 λ -parameter,

validate gives optimal

$MSE(\lambda, \text{polynomial complexity})$

Final model

Used on the independent test data.