

## Lecture September 4

- Bias-variance Tradeoff
- Resampling techniques
  - Bootstrap (Efron 1979)
  - Cross-validation
- Ridge Regression

### Bootstrap

resampling with replacement

$$\mathbf{Z} = (z_1, z_2, z_3, \dots, z_n) \quad \mathbf{Z} \subset A$$

$$P_n = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(z_i \in A)$$

### algorithm

(i) Draw a new sample

$$\mathbf{Z}_1^* = (z_1^*, z_2^*, \dots, z_n^*)$$

compute  $\hat{\theta}_1 = g(z_1^*, z_2^*, \dots, z_n^*)$

(ii) Repeat (i)  $B$  times,  
we have then

$$\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_B$$

(iii) compute standard deviation

$$STD = \sqrt{\frac{1}{B} \sum (\hat{\theta}_i - \bar{\theta})^2}$$

$$\bar{\theta} = \frac{1}{B} \sum_{j=1}^B \hat{\theta}_j$$

Discrete PDF :

$$E[x] = \sum_i x_i p(x_i) = \mu$$

$$\left( \int dx x p(x) \right)$$

sample mean  $\bar{\mu} = \frac{1}{n} \sum_i x_i$

$$p(x_i) = \frac{1}{n}$$

CV = cross-validation

LOOCV = Leave one out CV

$$X = \{x_1, x_2, \dots, x_n\}$$

$$Y = \{y_1, y_2, \dots, y_n\}$$

$$MSE = 0$$

for  $i=1, n$

$$X_{cv} = \{x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$$

$$Y_{cv} = \{y_1, y_2, \dots, y_{i-1}, y_{i+1}, \dots, y_n\}$$

$$x_{out} = x_i$$

$$y_{out} = \underset{\text{fit}}{\text{fit}}(X_{cv}, Y_{cv}, x_{out}) \quad \text{predict}$$

$$MSE = MSE + (y_i - \hat{y}_{i-1})^2$$

END FOR

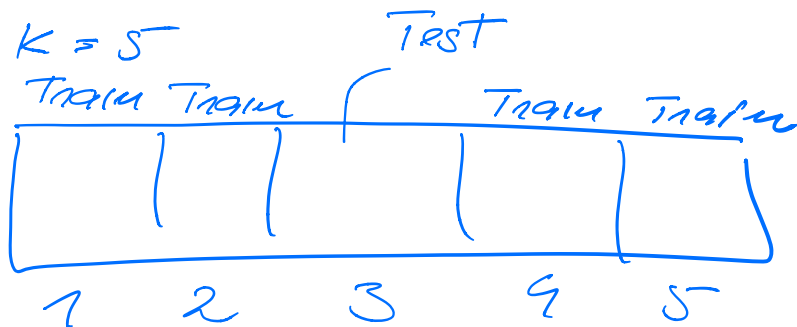
$$MSE = MSE/n$$

## K-Fold CV

Randomly partition data

K - equally sized sub  
samples

$$K = 5-10$$



The CV is then repeated  
K - times, with each of the  
K subsamplers used just  
once as test data,

$$K = 2$$

Model 1 :  $d_0, d_1$  (equal  
size)

Train on  $d_0$  and test

on  $d_1$  ( $\epsilon_{n1}$ )  
 Model 2 : Train on  $d_1$   
 and test on  $d_0$   
 ( $\epsilon_{n0}$ )

$$\text{Total error} = \text{err} = \frac{1}{2} (\text{err}_0 + \text{err}_1)$$

## Ridge Regression

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

if singular, no  $\hat{\beta}^{\text{OLS}}$

$$X^T X = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & & & \\ a_{m1} & \dots & \dots & a_{mn} \end{bmatrix}$$

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

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

$$(X^T X + \lambda \mathbb{I})$$

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

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

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

Corresponds to minimizing

$$\begin{aligned} C(\beta) &= \frac{1}{n} \sum_{i=0}^{n-1} (y_i - X_{i*} \beta)^2 \\ &\quad + \lambda \underbrace{\sum_{j=0}^{p-1} (\beta_j)^2}_{\beta^T \beta} \\ &= \frac{1}{n} \|y - X\beta\|_2^2 \\ &\quad + \lambda \|\beta\|_2^2 \end{aligned}$$

$\lambda =$  <sup>norm.</sup> penalty parameter /  
hyperparameter /  
regularization.

$$\lambda \geq 0$$

$L_2$ -type optimization.

$$\begin{aligned} \frac{\partial C(\beta)}{\partial \beta} &= 0 \quad - \frac{2}{n} X^T (y - X\beta) \\ &\quad + 2\lambda \beta = 0 \end{aligned}$$

$$\begin{aligned} \hat{\beta}^{\text{ridge}} &\Rightarrow \\ \hat{\beta} &= \underbrace{\left( X^T X + \lambda \underline{\underline{I}} \right)^{-1}}_{\text{not singular}} X^T y \end{aligned}$$

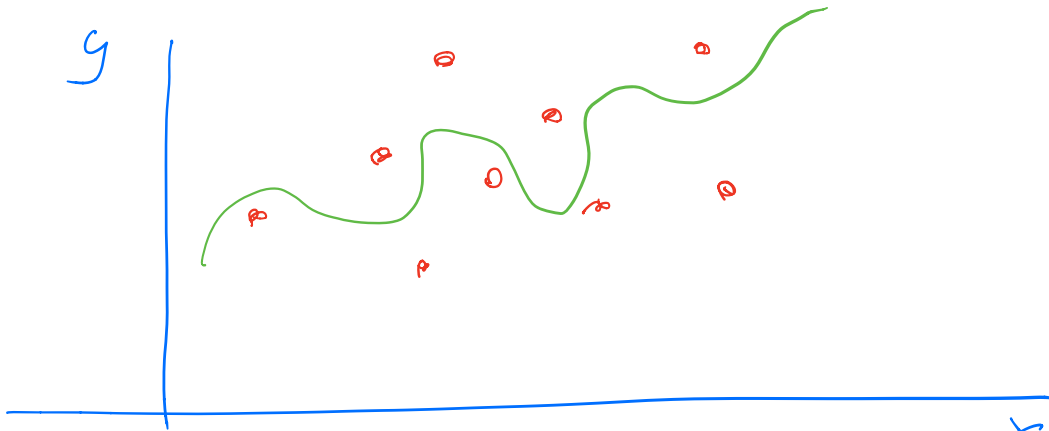
... - (ridge) - \)

(ridge)

$$MSE(\hat{\beta}^0(\lambda)) \leq MSE(\hat{\beta}^-(\lambda_{SO}))$$

What does this mean?

— Intuitive understanding



—  $\beta$ -values may show a large scatter in values large negative and large positive,

We could try to make these parameters small using Bayes's theorem:

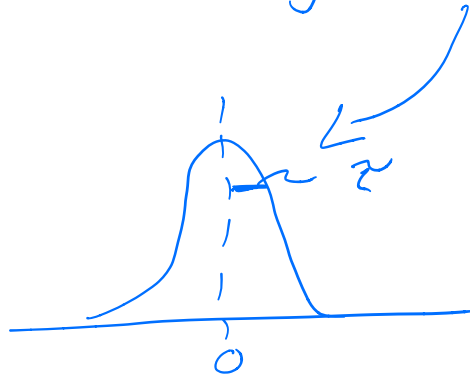
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

MLE

$$P(y_i | x_i; \beta) \propto \frac{e^{-\frac{(y_i - x_i \beta)^2}{2\sigma^2}}}{\sigma}$$

... P-1 ...

$$P(\beta) = \prod_{j=0}^p N(\beta_j | 0, \tau^2)$$



$$e^{-\beta_j^2 / 2\tau^2}$$

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

$$\begin{aligned} \text{MLE: } & P(y|x\beta) P(\beta) \\ &= \prod_{i=0}^{n-1} P(y_i | x_i \beta) \prod_{j=0}^{p-1} P(\beta_j) \end{aligned}$$

$$\begin{aligned} &= \sum_{i=0}^{n-1} \log P(y_i | x_i \beta) \\ &\quad + \sum_{j=0}^{p-1} \log P(\beta_j) + \text{const} \end{aligned}$$

$$\propto \sum_{i=0}^{n-1} (y_i - x_i \beta)^2 + \lambda \sum_{j=0}^{p-1} \beta_j^2$$

$$\lambda = \frac{1}{\tau^2} \geq 0$$

