

Lecture September 17

Resampling techniques

- Bootstrap (jackknife)
- cross validation

$$D = [x_1, x_2, \dots, x_n]$$

- (i) Draw a sample with replacement (randomly)

$$x_1^*, x_2^*, \dots, x_n^*$$

compute $\beta_{n,1}^* = g(x_1^*, x_2^*, \dots, x_n^*)$

$$\left[\mu_{n,1}^* = \frac{1}{n} \sum_{i=1}^n x_i^* \right]$$

- (ii) Repeat previous steps

B times

$$\beta_{n,1}^*, \beta_{n,2}^*, \dots, \beta_{n,B}^* \\ (\mu_{n,1}^* \dots)$$

- (iii) STD

$$S = \sqrt{\frac{1}{B} \sum_{j=1}^B (\beta_{n,j}^* - \bar{\beta})^2}$$

$$\bar{\beta} = \frac{1}{B} \sum_{j=1}^B \beta_{n,j}^*$$

First example

$$X \sim N(100, 15^2)$$

$$Z = \frac{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_m}{m} = \frac{m\bar{\mu}}{m} = \mu$$

$$X \sim N(\dots)$$

Central limit theorem

$$Z = \mu$$

$$\sigma_z^2 = \frac{\sigma^2}{m} \Rightarrow STD = \sigma/\sqrt{m}$$

in the code

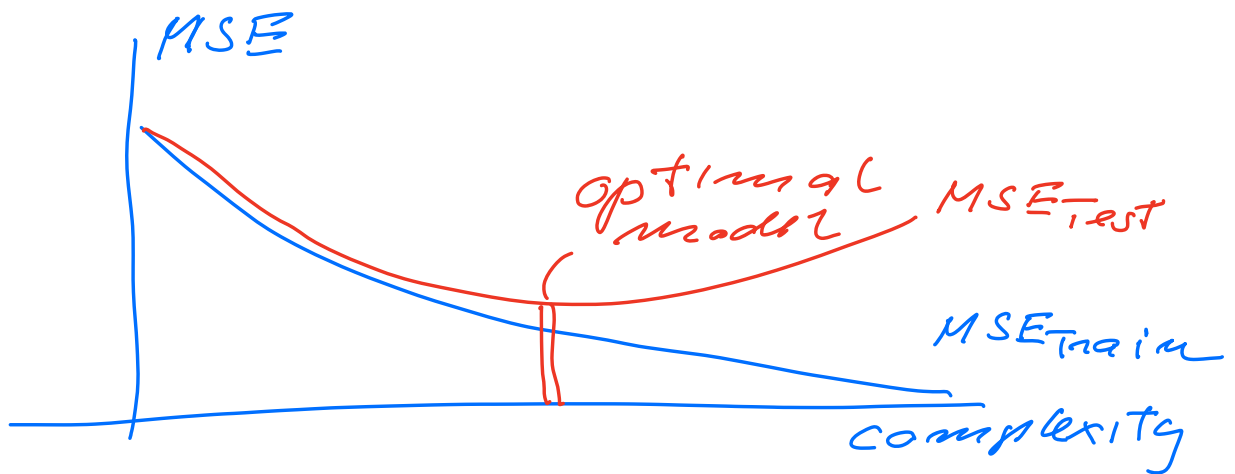
$$STD = \sigma/\sqrt{B}$$

Bias-variance Tradeoff

$$y_i = f(x_i) + \varepsilon_i$$

$$f(x_i) \approx \tilde{y}_i \quad (\text{Model } x_i \beta)$$

$$\begin{aligned} \text{MSE} &= E[(y - X\beta)^2] \\ &= \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 \end{aligned}$$



Bootstrap to improve
(or make reliable prediction
of MSE_{test})

Bias-variance tradeoff:

$$n-1 \quad \sim \quad \frac{1}{n}$$

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

$$\left(\begin{aligned} E[\tilde{y}_i] &= X_i^T \beta = \tilde{y}_i \\ \text{var}[\tilde{y}_i] &= \text{var}[y_i] = \\ \text{var}[\varepsilon_i] &= \sigma^2 \\ E[\hat{\beta}] &= \hat{\beta} \quad \text{var}[\hat{\beta}]_{OLS} = \sigma^2 (X^T X)^{-1} \end{aligned} \right)$$

$$\text{var}[\tilde{y}_i] = \text{var}[y_i] =$$

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

$$E[\hat{\beta}] = \hat{\beta} \quad \text{var}[\hat{\beta}]_{OLS} = \sigma^2 (X^T X)^{-1}$$

$$= \text{Bias} + \text{variance} + \sigma^2$$

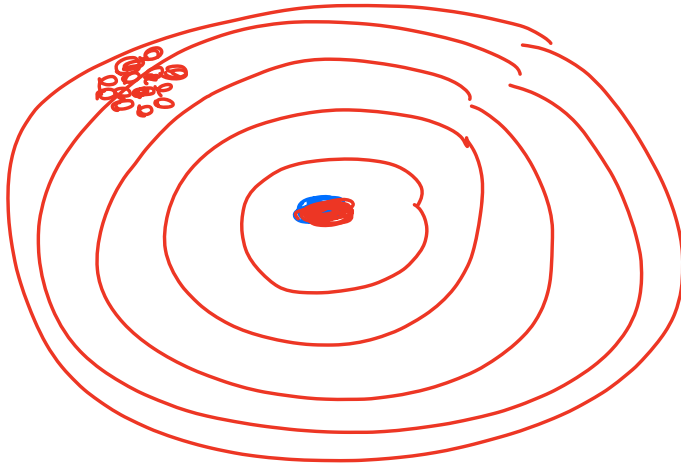
$$\frac{1}{n} \sum_{i=0}^{n-1} (y_i - E[\tilde{y}])^2$$

$$\frac{1}{n} \sum_{i=0}^{n-1} (\tilde{y}_i - E[\tilde{y}])^2$$

$$E[\tilde{y}] = \frac{1}{n} \sum_{i=0}^{n-1} \tilde{y}_i$$

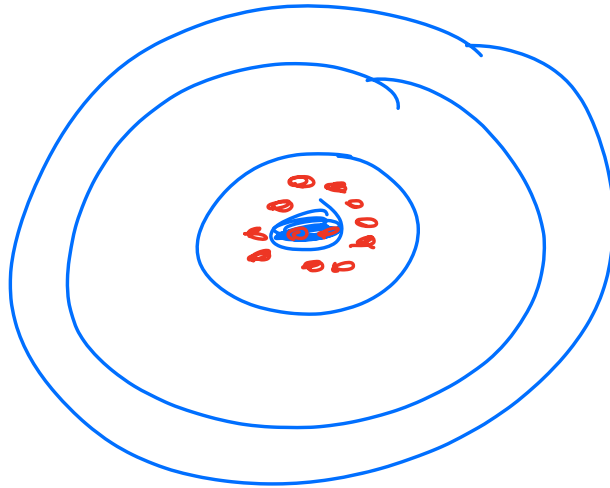
Low variance

High
Bias



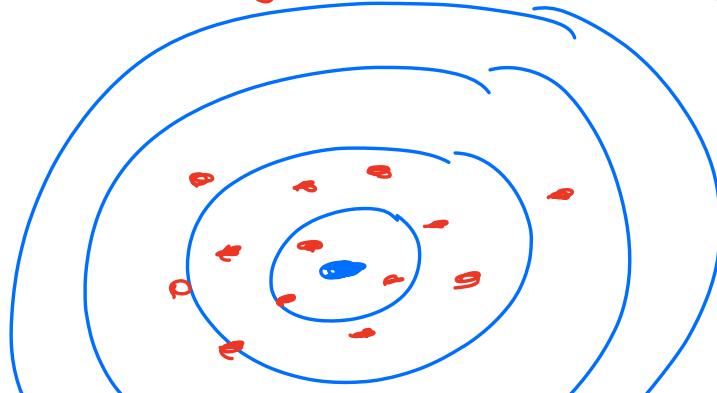
Low variance

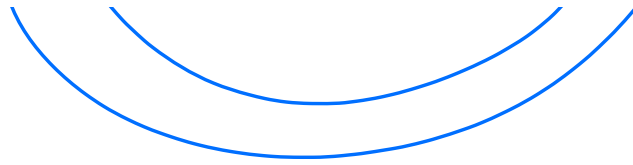
Low
Bias



High variance

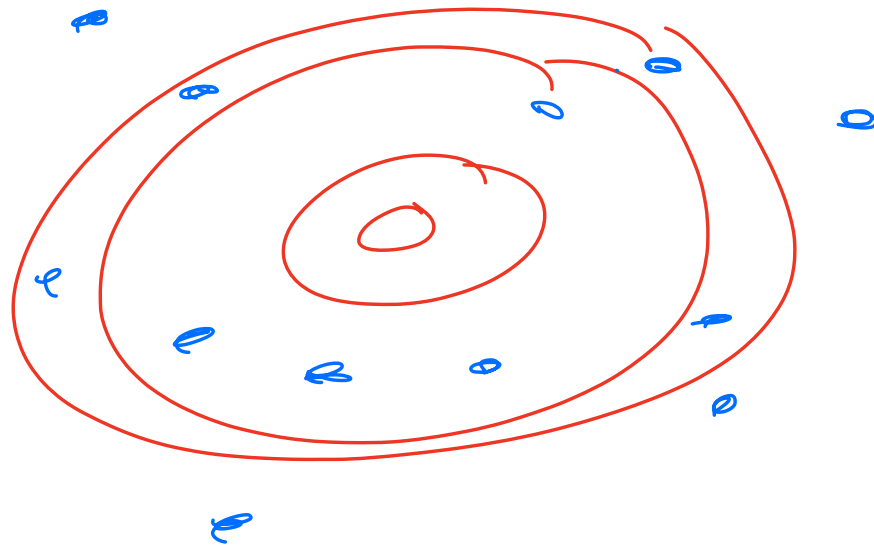
Low
Bias





High variance

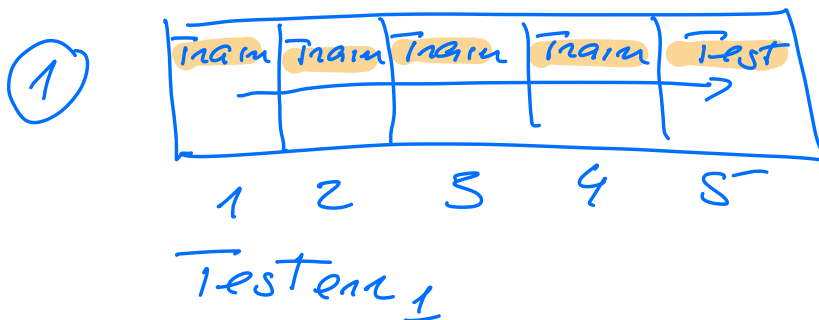
High Bias



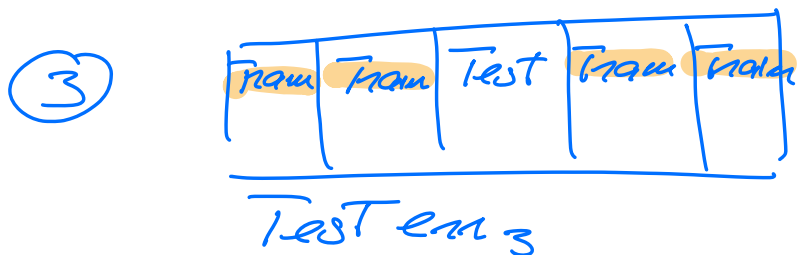
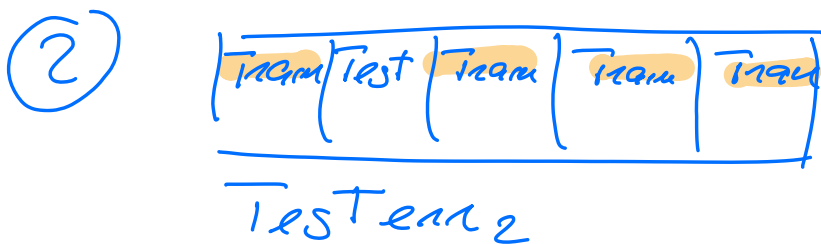
Cross-validation (CV)

k-fold CV

$k = 5$



Data set

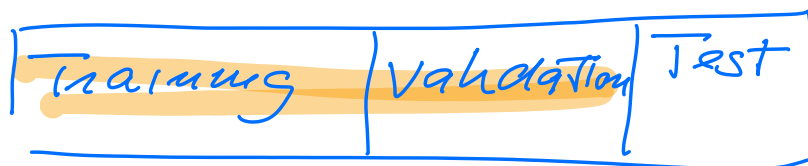


⑤

$$Err = \frac{1}{K} \sum_{i=1}^K TestErr_i$$

$K = 5 - 10$

splitting of data



$k = n$ (number of data)
LOCV

$$x = [x_1, x_2, \dots, x_n]$$

$$y = [y_1, y_2, \dots, y_n]$$

for $i = 1, n$

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

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

$$x_{out} = x_i$$

$$\tilde{y}_{out} = f_1^+(x_{cv}, y_{cv}, x_{out})$$