

Lecture FYS-STK,
September 14,
2023

Resampling (Bootstrap, cross-validation)

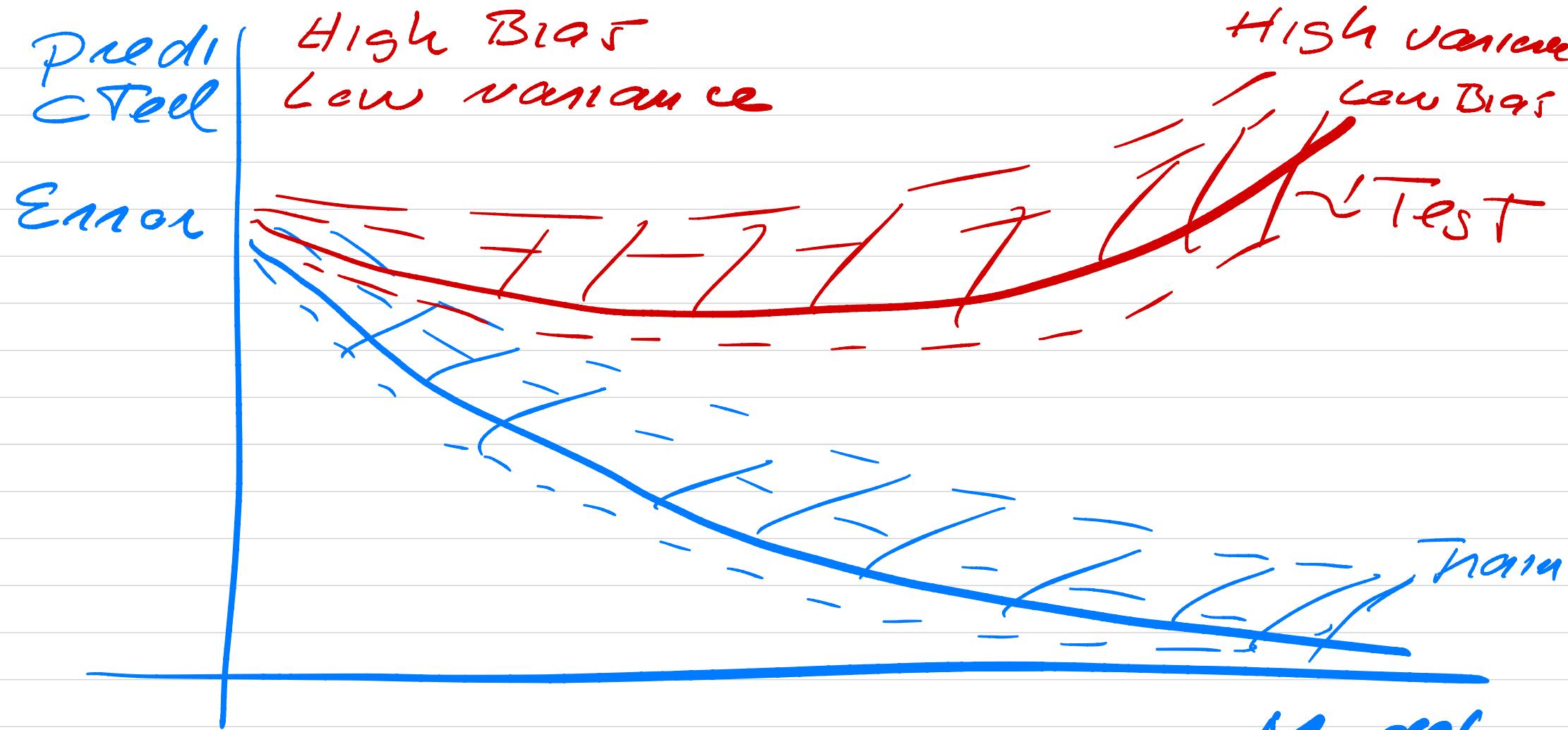
We have defined functions to assess model, MSE

$$L(y_i, \hat{y}_i)$$

\nwarrow model

with given training data T_i and test data \tilde{T}_i

$$E_{\text{test}} = L(y_i, \hat{y}_i; T_i \tilde{T}_i)$$



$$\text{Err} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, \hat{y}_i; \theta_i, \gamma_i)$$

Bias-variance tradeoff for MSE

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

$$= [E[(y - \hat{y})^2]]$$

$$= [E[y^2]] - 2[E[yy]] + [E[\hat{y}^2]]$$

$$[E[\hat{y}^2]] = \text{var}[\hat{y}] + (E[\hat{y}])^2$$

$$\begin{aligned} E[y^2] &= E[(f+\varepsilon)^2] = \\ &= E[f^2] + 2[E[f\varepsilon]] + E[\varepsilon^2] = \sigma^2 \end{aligned}$$

$$= f^2 + 2f \cdot 0 + \sigma^2$$

$$E[\hat{y}\hat{y}] = E[(f+\varepsilon)\hat{y}]$$

$$= E[f\hat{y}] + E[\varepsilon\hat{y}]$$

$$= fE[\hat{y}] + E[\hat{y}]E[\varepsilon]$$

Collecting ≈ 0

$$\begin{aligned} & f^2 - 2fE[\hat{y}] + \sigma^2 + \text{Var}[\hat{y}] \\ & + (E[\hat{y}])^2 \end{aligned}$$

$$= E[(f - E[\tilde{g}])^2] + \sigma^2$$

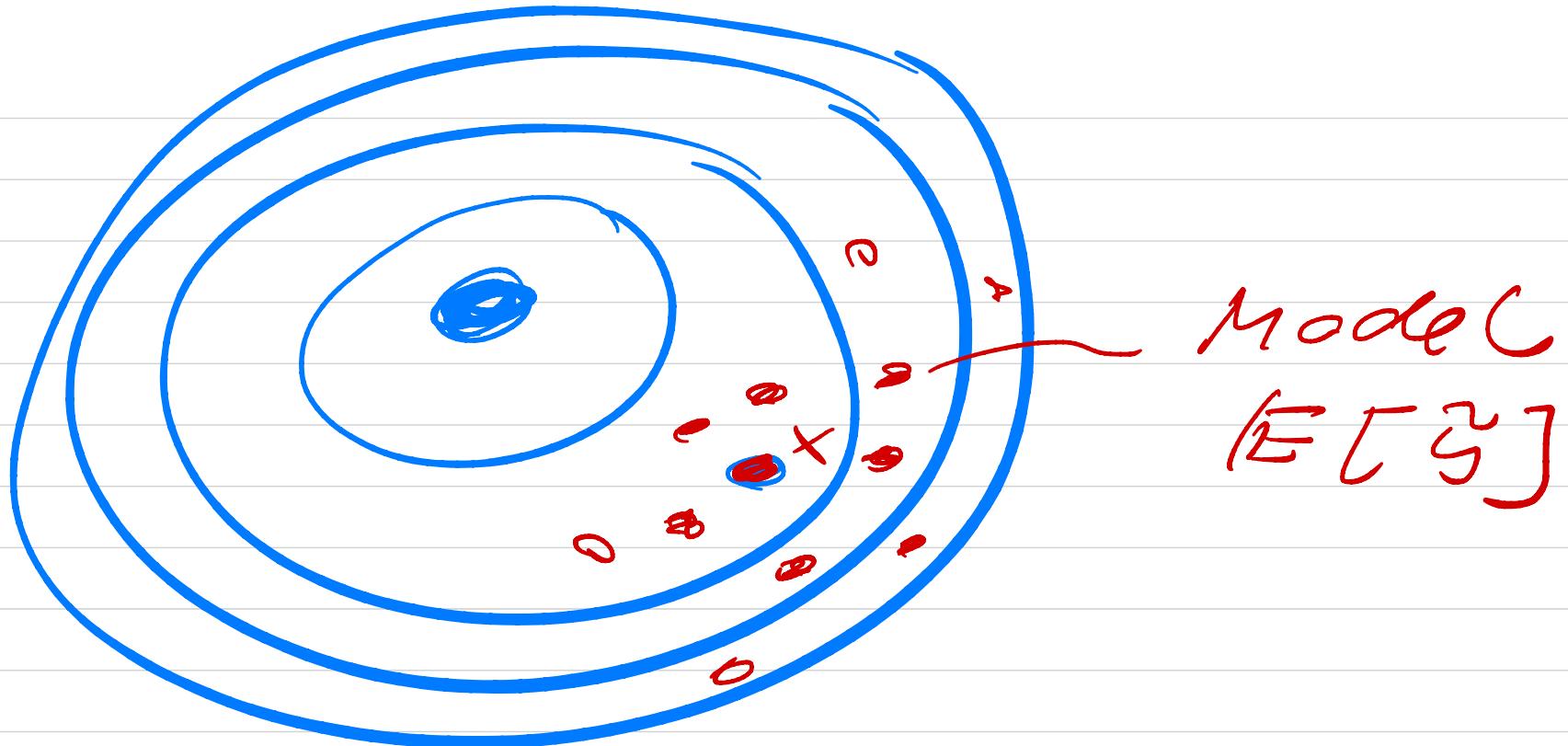
+ $\text{var}[\tilde{g}]$

Bias

variance

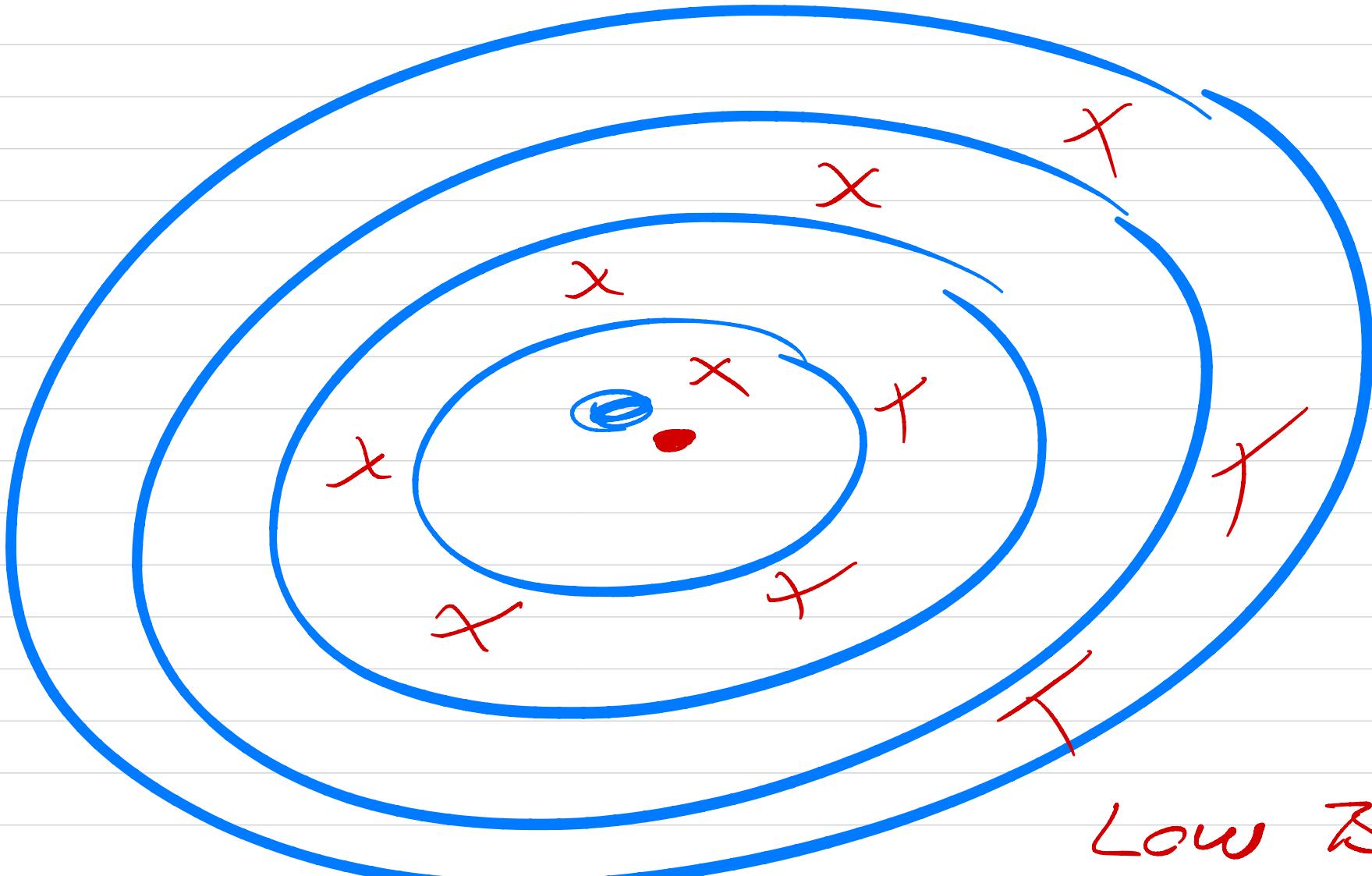
$$\text{Bias} \approx \frac{1}{n} \sum_{i=0}^{n-1} (y_i - E[\tilde{g}])^2$$

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

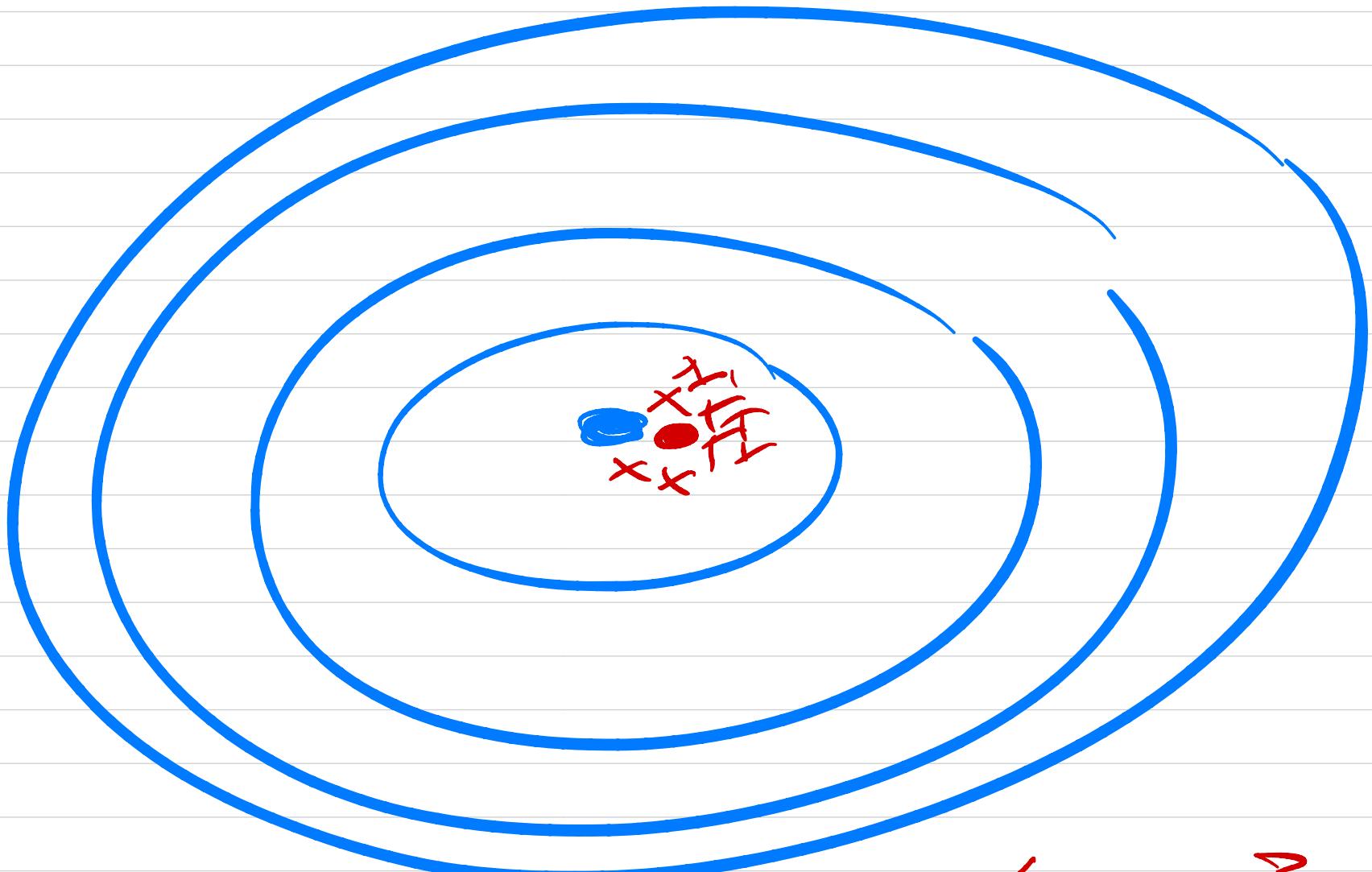


High Bias
High variance





Low Bias
High variance



Low Bias
low
variance

Bootstrap

Original data set

$$D: \{z_1, z_2, \dots, z_N\}$$

i.i.d.

(i) Draw a Bootstrap sample

$$D_1^* = \{z_1^*, z_2^*, \dots, z_n^*\}$$

with replacement

Compute \hat{e}_1^* (mean value)

(ii) Repeat B times
getting estimator

$$\hat{\theta}_1^* \quad \hat{\theta}_2^* \quad \dots \quad \hat{\theta}_B^*$$

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

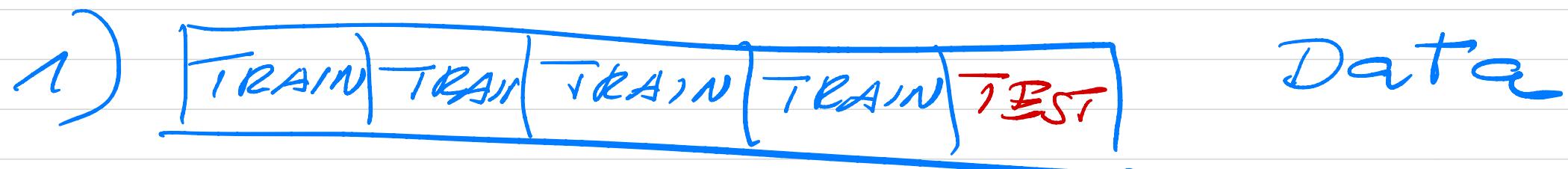
(iii) can compute variance

$$\sigma_\theta^2 = \frac{1}{B} \sum_{j=1}^B (\hat{\theta}_j^* - \bar{\theta})^2$$

Resampling with cross-validation

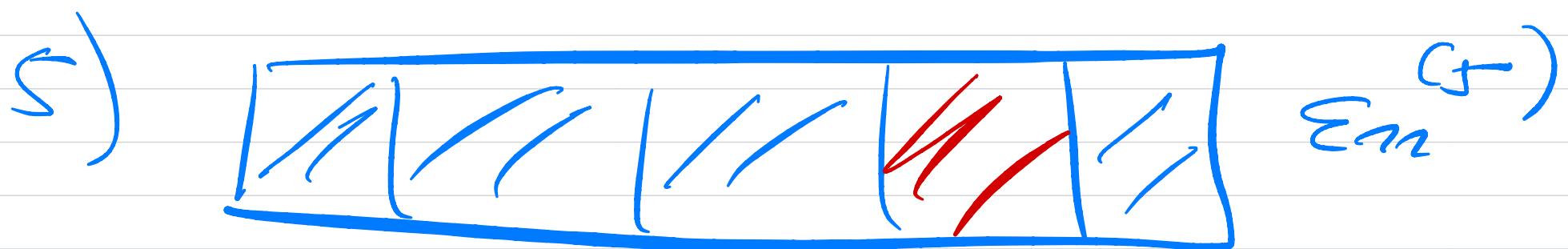
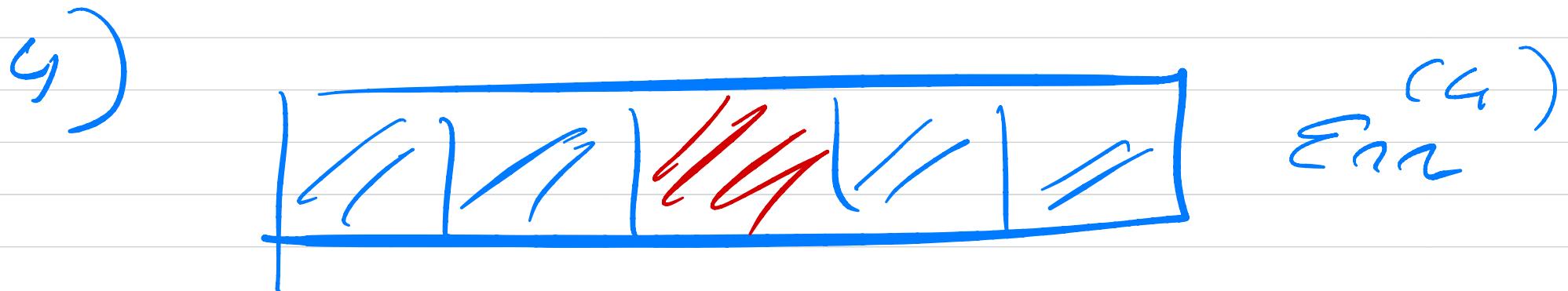
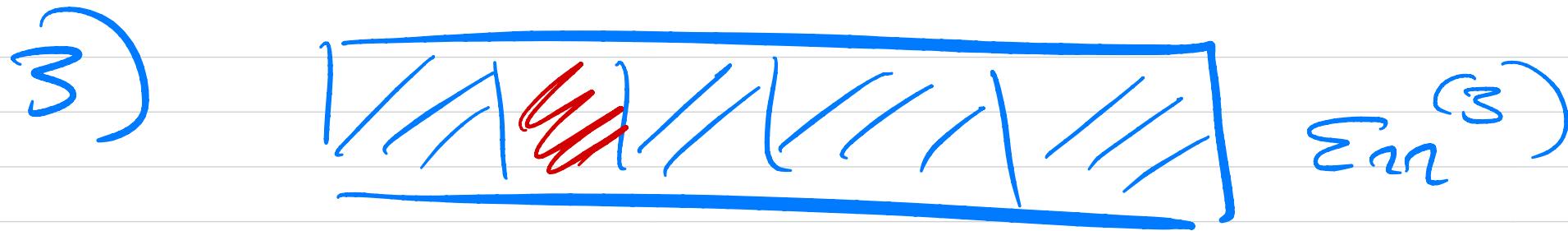
K-Folds : splitting data
in K-subsets

$$\underline{K=5}$$



$$\epsilon_{\text{nnTest}} = \sum_n^{(1)}$$





$$\Sigma m = \frac{1}{5} \sum_{i=1}^5 \Sigma m^{(i)}$$

Typical values of $K \approx 5-10$

But need to test that MSE

stabilizes -