



Machine Learning 1

Lecture 4.2 - Supervised Learning
Bias Variance Decomposition

Erik Bekkers

(Bishop 1.5.5, 3.2)



Expected Loss for Regression

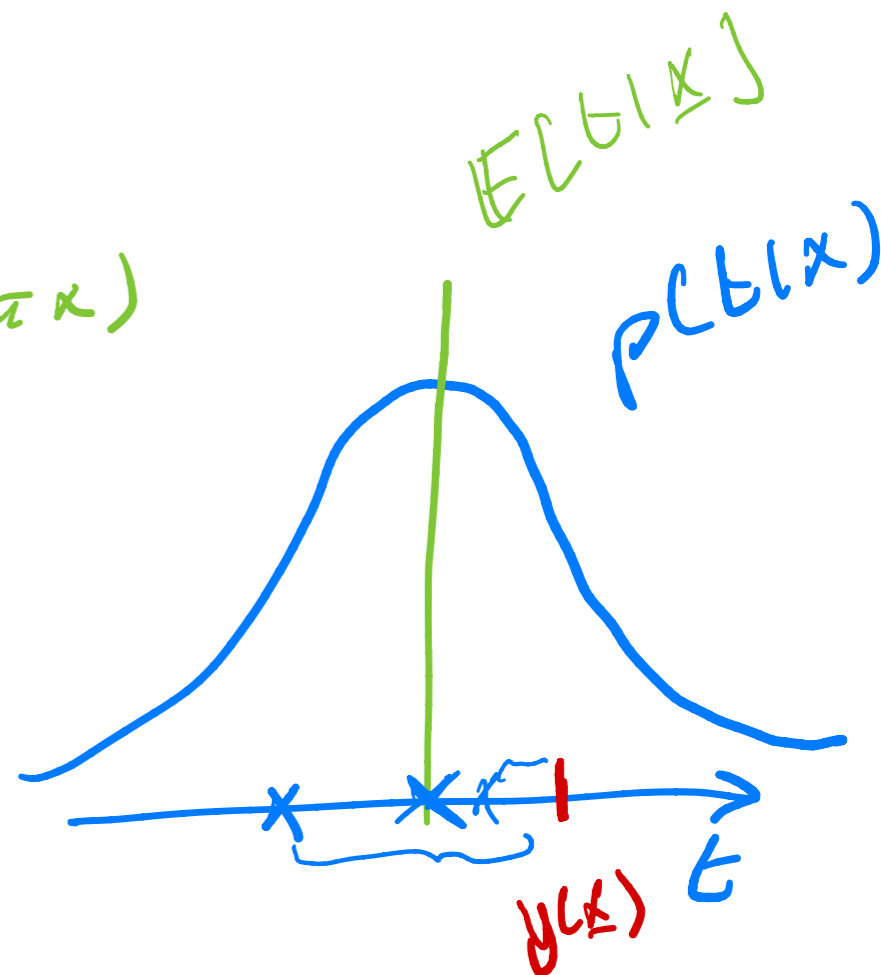
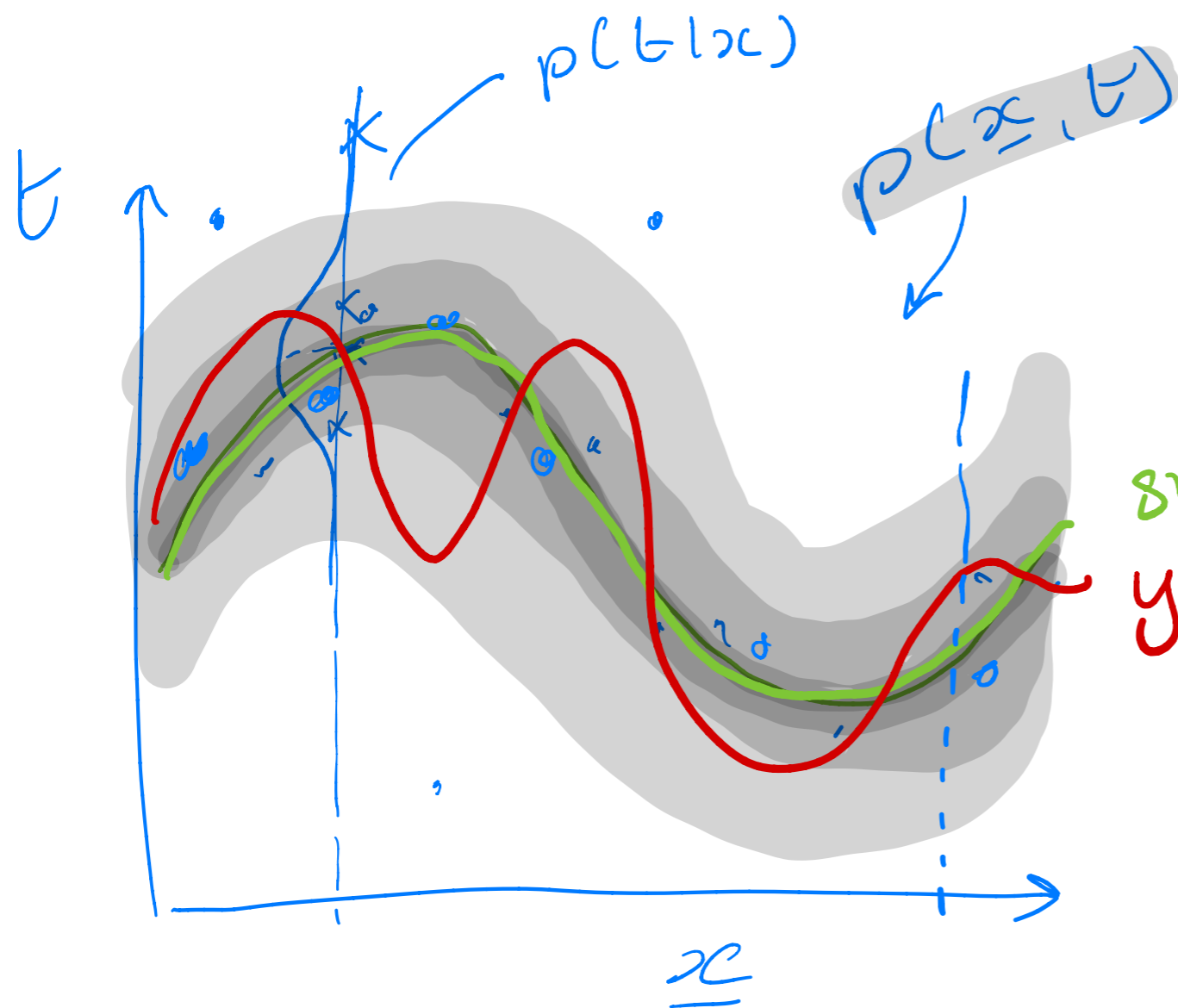
Frequentist viewpoint of model complexity

for a given $(\underline{x}, t) \sim p(\underline{x}, t)$

- ▶ Regression loss function: $L(t, y(\mathbf{x})) = (t - y(\underline{x}))^2$
- ▶ Expected loss: $\mathbb{E}[L(t, y(\mathbf{x}))] = \iint \underbrace{(t - y(\underline{x}))^2}_{\text{red underline}} \underbrace{p(\underline{x}, t)}_{\text{green underline}} \underbrace{d\underline{x} dt}_{\text{blue underline}}$
- ▶ Optimal $y(\mathbf{x})$ minimizes $\mathbb{E}[L(t, y(\mathbf{x}))]$

$$y(\underline{x}) = \mathbb{E}[t | \underline{x}]$$

$$t = \sin(2\pi x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$



For fixed \underline{x} , the expected loss

minimum w.r.t. $y(\underline{x})$

$$\frac{\partial}{\partial y(\underline{x})} \dots = 0 \Rightarrow \int y(\underline{x}) p(t|\underline{x}) dt = \int t p(t|\underline{x}) dt$$

regression function $y(\underline{x}) = E[t|\underline{x}]$

Expected Loss for Regression

- ▶ Decomposition of expected loss:

Verify Bishop 10.5

$$\mathbb{E}[L] = \int \int \underbrace{(y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])}_{\text{bias}} + \underbrace{\mathbb{E}[t|\mathbf{x}] - t}_{\text{due to intrinsic noise}})^2 p(\mathbf{x}, t) dt d\mathbf{x}$$
$$= \int \underbrace{(y(\mathbf{x}) - E[t|\mathbf{x}])^2}_{\text{bias}} p(\mathbf{x}) d\mathbf{x} + \int \underbrace{\text{var}[t|\mathbf{x}] p(\mathbf{x})}_{\text{due to intrinsic noise}} d\mathbf{x}$$

Minimizing the Expected Loss

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- ▶ Optimal solution is $y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$ (unknown)
- ▶ Only finite dataset observed: $\{(x_1, t_1), \dots, (x_n, t_n)\} = D$
.. .. = D_2
- ▶ Frequentist approach: estimate $y_D(\mathbf{x}) = y(\mathbf{x}, \mathbf{w}^*)$ based on dataset D
- ▶ Estimate performance of learning algorithm by averaging the expected loss over learned $y_D(\mathbf{x})$ for different datasets D

$$\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2]$$

Bias-Variance Decomposition

$$\mathbb{E}[\mathbb{E}_D[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2] p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- ▶ Bias-Variance decomposition:

$$\mathbb{E}_D[\{y_D(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2] = \mathbb{E}_D[\underbrace{\{y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})]\}}_{\text{variance}} + \underbrace{\{\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t|\mathbf{x}]\}}_{\text{bias}}]^2 =$$

$$= \mathbb{E}_D[\underbrace{(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2}_{\text{variance}} + \underbrace{(\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t|\mathbf{x}])^2}_{\text{bias}}]$$

$$= \mathbb{E}_D[\underbrace{\{y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})]\}}_{\text{variance}} \underbrace{\{\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t|\mathbf{x}]\}}_{\text{bias}}]$$

- ▶ Expected loss decomposition: $\mathbb{E}[\mathbb{E}_D[L]] = (\text{bias})^2 + \text{variance} + \text{noise}$

$$\text{bias}^2 = \int (\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int (\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2]) p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

Bias-Variance Decomposition: Example

- Generate L datasets of N points: $L=100$, $N=25$

$$x \sim U(0, 1)$$

$$t = \sin(2\pi x) + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \alpha^{-1})$$

$$\mathbb{E}[t|x] = \sin(2\pi x)$$

- L predictions with 24 Gaussian basis functions

$$y^{(l)}(x) = (\mathbf{w}^{(l)})^T \phi(x)$$

$$E_D = \frac{1}{2} \sum_{i=1}^N \{t_n - \mathbf{w}^T \phi(x)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\mathbb{E}_D[y_D(x)] = \bar{y}(x)$$

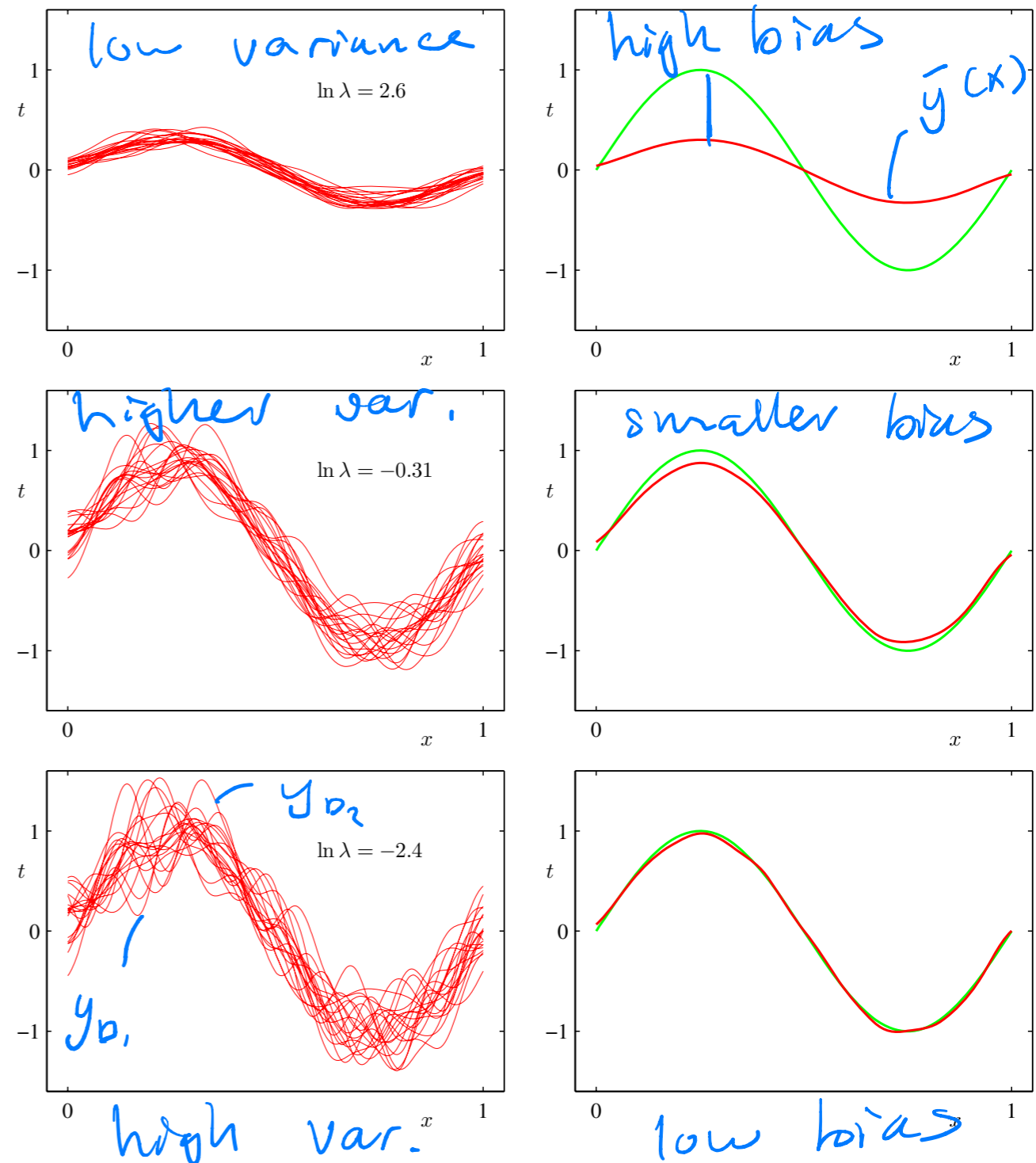


Figure: bias-variance decomposition (Bishop 3.5)

Bias-Variance Decomposition: Example

Estimate the bias and variance:

$$\begin{aligned} \text{bias}^2 &= \int \left\{ \mathbb{E}_D[y_D(x)] - \mathbb{E}[t|x] \right\}^2 p(x) dx \\ &= \frac{1}{N} \sum_{n=1}^N \left(\bar{y}(x_n) - \mathbb{E}[t|x_n] \right)^2 \end{aligned}$$

Handwritten notes: $\bar{y}(x)$ is written above the first term, $\mathbb{E}[t|x]$ is written above the second term, and $\{x_1, x_2, \dots, x_N\}$ is written in green with an arrow pointing to the sum.

$$\mathbb{E}_D[y_D(x)] = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x) = \bar{y}(x)$$

$$\begin{aligned} \text{variance} &= \int \mathbb{E}_D \left[\left\{ y_D(x) - \mathbb{E}_D[y_D(x)] \right\}^2 \right] p(x) dx \\ &= \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \left(y^{(l)}(x_n) - \bar{y}(x_n) \right)^2 \end{aligned}$$

Bias-Variance Decomposition: Example

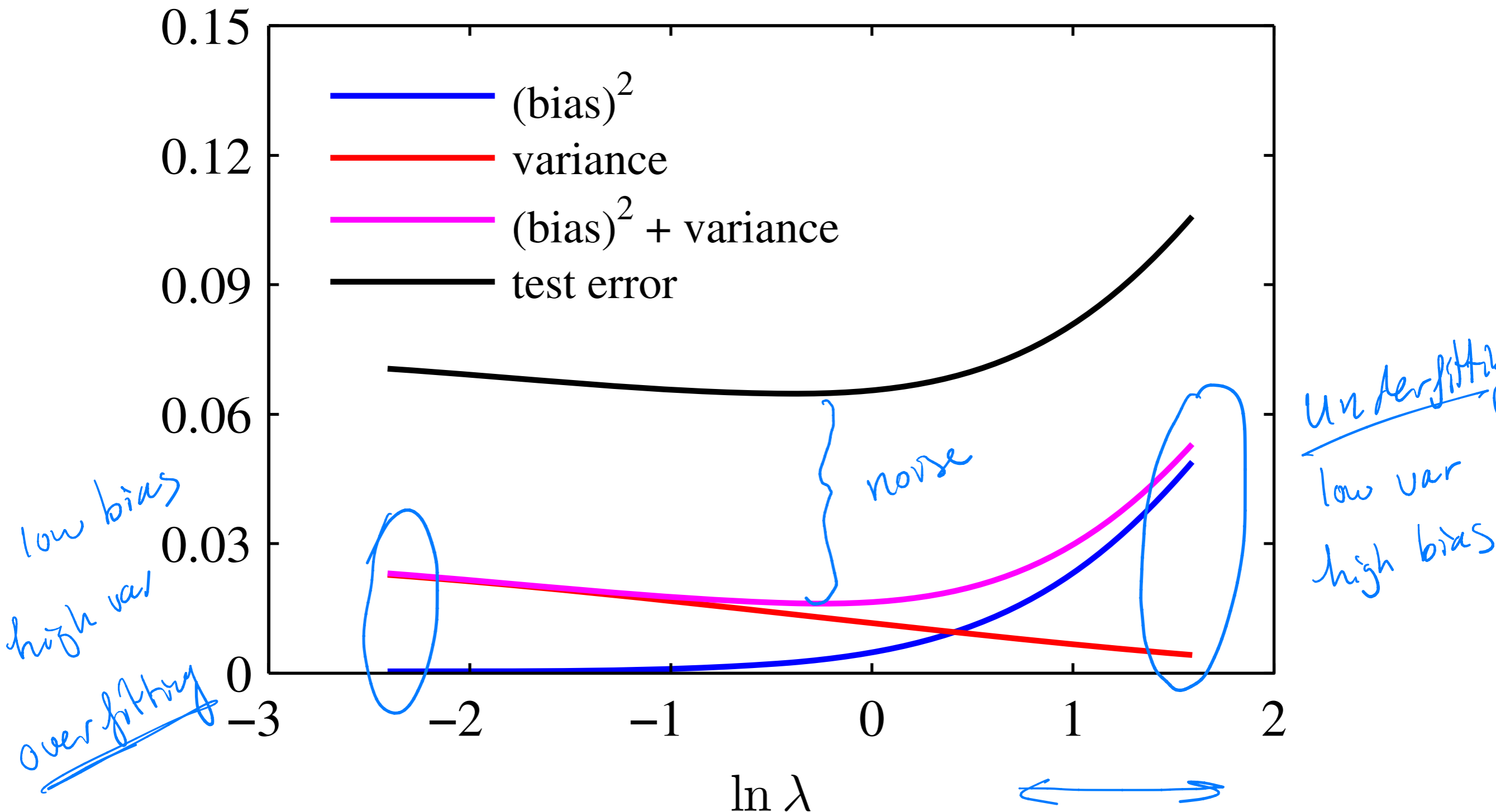


Figure: bias-variance decomposition (Bishop 3.6)

trade off
var / bias

Bias-Variance decomposition

- ▶ In practice we don't want to split our dataset into L datasets to determine the best model complexity (best value of λ)
- ▶ Better to keep large dataset,
 - ▶ Less overfitting.
 - ▶ Different optimal model complexity!
- ▶ Bayesian regression!