

Bias-Variance Trade-off, Regularization

Machine Learning

Hamid R Rabiee – Zahra Dehghanian
Spring 2025



Sharif University
of Technology

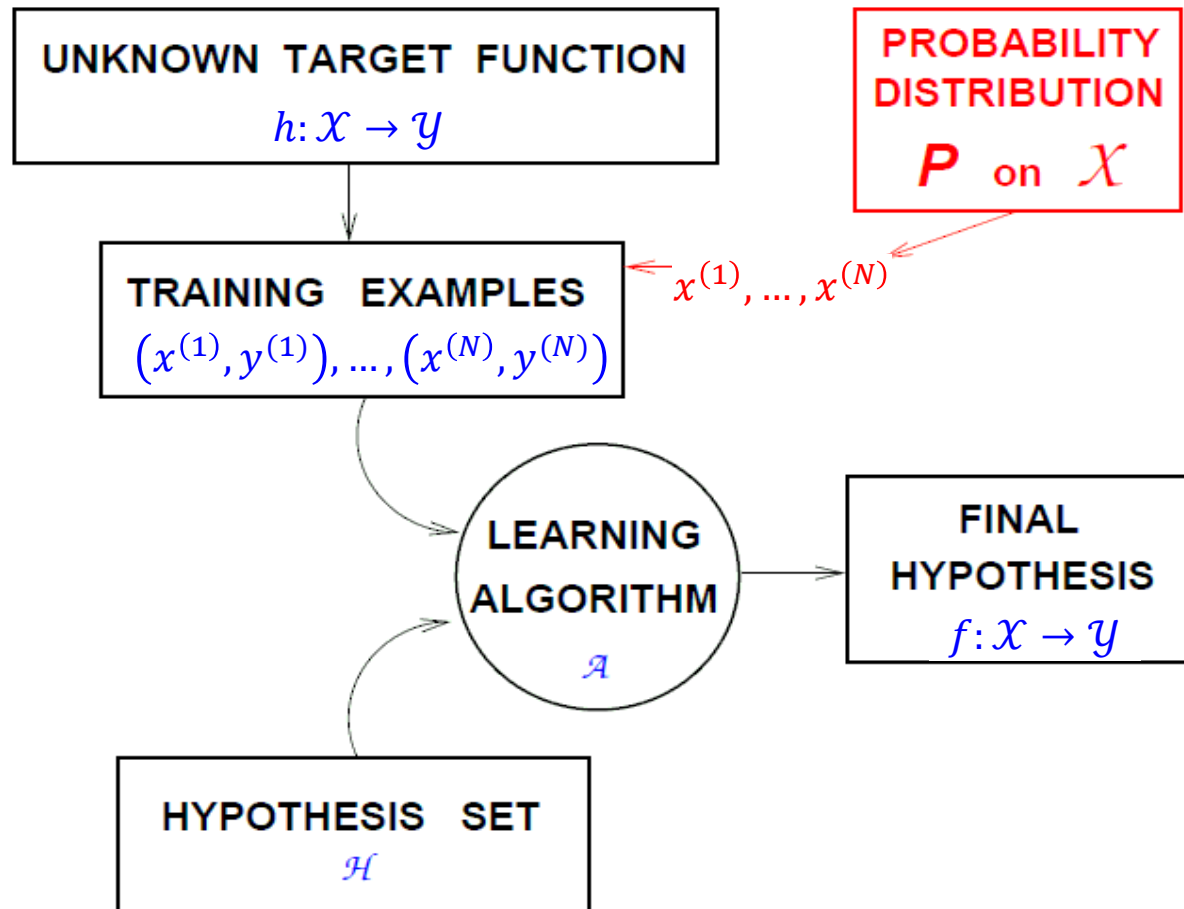
Model complexity: Bias-variance trade-off

- ▶ Least squares can lead to severe over-fitting if complex models are trained using data sets of limited size.
- ▶ A frequentist viewpoint of the model complexity issue, known as the *bias-variance trade-off*.

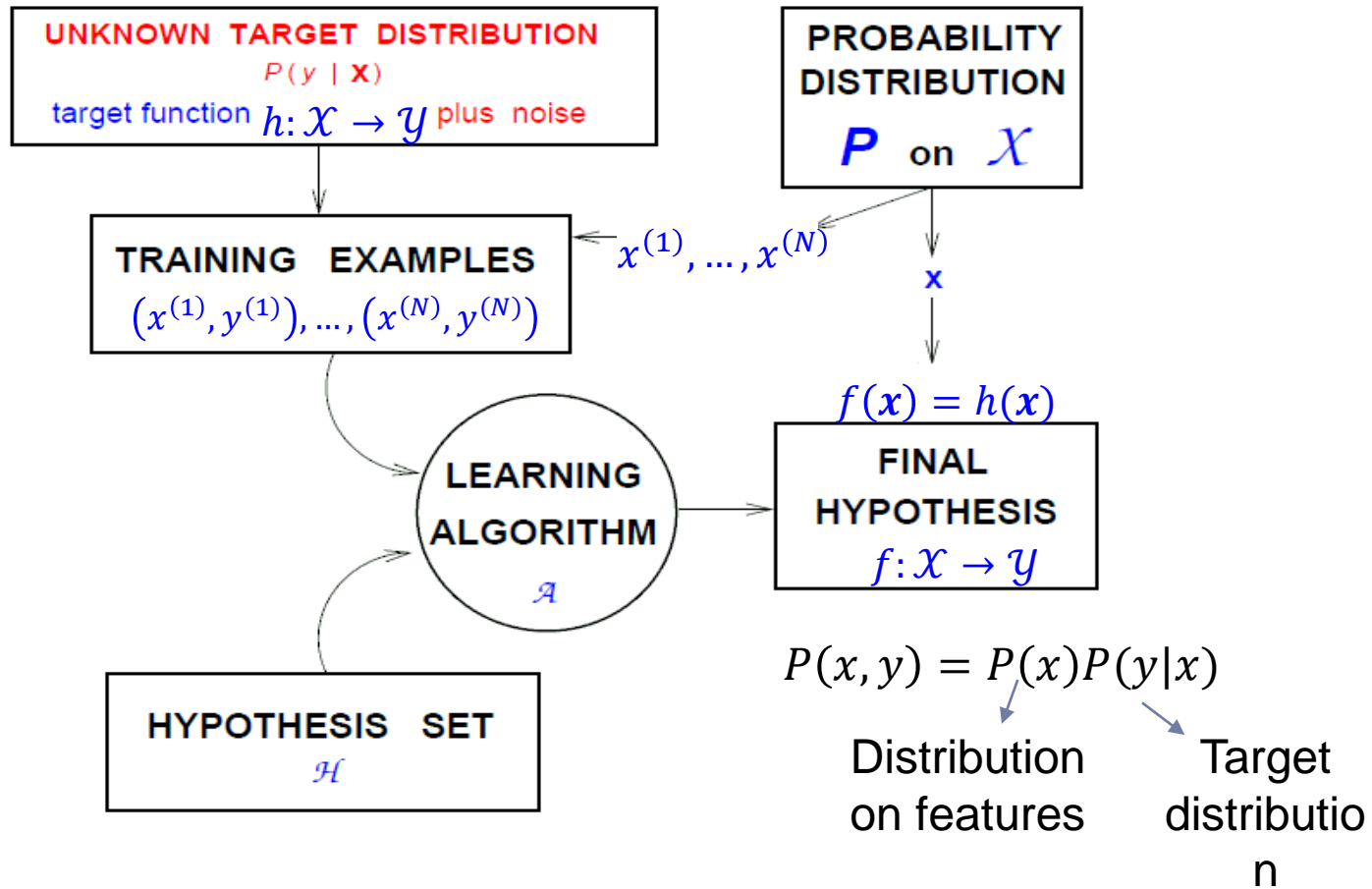
Formal discussion on bias, variance, and noise

- ▶ Best unrestricted regression function
- ▶ Noise
- ▶ Bias and variance

The learning diagram: deterministic target



The learning diagram including noisy target



[Y.S. Abou Mostafa, 2012]

Best unrestricted regression function

- ▶ If we know the joint distribution $P(\mathbf{x}, y)$ and no constraints on the regression function?
 - ▶ cost function: mean squared error

$$h^* = \operatorname{argmin}_{h: \mathbb{R}^d \rightarrow \mathbb{R}} \mathbb{E}_{\mathbf{x}, y} \left[(y - h(\mathbf{x}))^2 \right]$$

$$h^*(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y]$$

Best unrestricted regression function: Proof

$$\mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right] = \iint (y - h(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

Best unrestricted regression function: Proof

$$\mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right] = \iint (y - h(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

- For each \mathbf{x} , separately minimize loss since $h(\mathbf{x})$ can be chosen independently for each different \mathbf{x} :

$$\frac{\delta \mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right]}{\delta h(\mathbf{x})} = - \int 2(y - h(\mathbf{x})) p(\mathbf{x}, y) dy = 0$$

Best unrestricted regression function: Proof

$$\mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right] = \iint (y - h(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

- For each \mathbf{x} , separately minimize loss since $h(\mathbf{x})$ can be chosen independently for each different \mathbf{x} :

$$\begin{aligned} \frac{\delta \mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right]}{\delta h(\mathbf{x})} &= - \int 2(y - h(\mathbf{x})) p(\mathbf{x}, y) dy = 0 \\ \Rightarrow h(\mathbf{x}) &= \frac{\int y p(\mathbf{x}, y) dy}{\int p(\mathbf{x}, y) dy} = \frac{\int y p(\mathbf{x}, y) dy}{p(\mathbf{x})} = \int y p(y|\mathbf{x}) dy = \mathbb{E}_{y|\mathbf{x}} [y] \end{aligned}$$

Best unrestricted regression function: Proof

$$\mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right] = \iint (y - h(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

- For each \mathbf{x} , separately minimize loss since $h(\mathbf{x})$ can be chosen independently for each different \mathbf{x} :

$$\frac{\delta \mathbb{E}_{\mathbf{x},y} \left[(y - h(\mathbf{x}))^2 \right]}{\delta h(\mathbf{x})} = - \int 2(y - h(\mathbf{x})) p(\mathbf{x}, y) dy = 0$$

$$\Rightarrow h(\mathbf{x}) = \frac{\int y p(\mathbf{x}, y) dy}{\int p(\mathbf{x}, y) dy} = \frac{\int y p(\mathbf{x}, y) dy}{p(\mathbf{x})} = \int y p(y|\mathbf{x}) dy = \mathbb{E}_{y|\mathbf{x}} [y]$$

$$\Rightarrow h^*(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}} [y]$$

Error decomposition

$$(x, y) \sim P$$

$h(x)$: minimizes the expected loss

$$E_{true}(f_{\mathcal{D}}(x)) = \mathbb{E}_{x,y}[(f_{\mathcal{D}}(x) - y)^2] \quad \text{Expected loss}$$

Error decomposition

$$(\mathbf{x}, y) \sim P$$

$h(\mathbf{x})$: minimizes the expected loss

$$E_{true}(f_{\mathcal{D}}(\mathbf{x})) = \mathbb{E}_{\mathbf{x}, y}[(f_{\mathcal{D}}(\mathbf{x}) - y)^2] \quad \text{Expected loss}$$

$$= \mathbb{E}_{\mathbf{x}, y} [(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) + h(\mathbf{x}) - y)^2]$$

Error decomposition

$$(x, y) \sim P$$

$h(x)$: minimizes the expected loss

$$E_{true}(f_{\mathcal{D}}(x)) = \mathbb{E}_{x,y}[(f_{\mathcal{D}}(x) - y)^2]$$

Expected loss

$$= \mathbb{E}_{x,y} [(f_{\mathcal{D}}(x) - h(x) + h(x) - y)^2]$$

$$\begin{aligned} &= \mathbb{E}_x \left[(f_{\mathcal{D}}(x) - h(x))^2 \right] + \mathbb{E}_{x,y} [(h(x) - y)^2] \\ &\quad + 2\mathbb{E}_{x,y} [(f_{\mathcal{D}}(x) - h(x))(h(x) - y)] \end{aligned}$$



Error decomposition

$(\mathbf{x}, y) \sim P$

$h(\mathbf{x})$: minimizes the expected loss

$$E_{true}(f_{\mathcal{D}}(\mathbf{x})) = \mathbb{E}_{\mathbf{x}, y}[(f_{\mathcal{D}}(\mathbf{x}) - y)^2]$$

Expected loss

$$= \mathbb{E}_{\mathbf{x}, y} [(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) + h(\mathbf{x}) - y)^2]$$

$$= \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathbf{x}, y} [(h(\mathbf{x}) - y)^2] \\ + 2\mathbb{E}_{\mathbf{x}, y} [(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))(h(\mathbf{x}) - y)]$$

$$\underbrace{\hspace{15em}}_{\mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x})) \mathbb{E}_{y|x} [(h(\mathbf{x}) - y)] \right]}$$

Error decomposition

$$(x, y) \sim P$$

$h(x)$: minimizes the expected loss

$$E_{true}(f_{\mathcal{D}}(x)) = \mathbb{E}_{x,y}[(f_{\mathcal{D}}(x) - y)^2]$$

Expected loss

$$= \mathbb{E}_{x,y} [(f_{\mathcal{D}}(x) - h(x) + h(x) - y)^2]$$

$$= \mathbb{E}_x \left[(f_{\mathcal{D}}(x) - h(x))^2 \right] + \mathbb{E}_{x,y} [(h(x) - y)^2] \\ + 2\mathbb{E}_{x,y} [(f_{\mathcal{D}}(x) - h(x))(h(x) - y)]$$



$$\mathbb{E}_x \left[(f_{\mathcal{D}}(x) - h(x)) \underbrace{\mathbb{E}_{y|x}[(h(x) - y)]}_{0} \right]$$

0



Error decomposition

$$(x, y) \sim P$$

$h(x)$: minimizes the expected loss

$$\begin{aligned} E_{true}(f_{\mathcal{D}}(x)) &= \mathbb{E}_{x,y}[(f_{\mathcal{D}}(x) - y)^2] \\ &= \mathbb{E}_{x,y}[(f_{\mathcal{D}}(x) - h(x) + h(x) - y)^2] \\ &= \mathbb{E}_x \left[(f_{\mathcal{D}}(x) - h(x))^2 \right] + \underbrace{\mathbb{E}_{x,y}[(h(x) - y)^2]}_{\text{noise}} \\ &\quad + 0 \end{aligned}$$

- ▶ Noise shows the irreducible minimum value of the loss function

Expectation of true error

$$\begin{aligned} E_{true}(f_{\mathcal{D}}(\mathbf{x})) &= \mathbb{E}_{\mathbf{x},y}[(f_{\mathcal{D}}(\mathbf{x}) - y)^2] \\ &= \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] + noise \end{aligned}$$

Expectation of true error

$$\begin{aligned} E_{true}(f_{\mathcal{D}}(\mathbf{x})) &= \mathbb{E}_{\mathbf{x},y}[(f_{\mathcal{D}}(\mathbf{x}) - y)^2] \\ &= \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] + noise \end{aligned}$$

$$\begin{aligned} &\mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] \right] \end{aligned}$$

We now want to focus on $\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right]$.

The average hypothesis

$$\bar{f}(\mathbf{x}) \equiv E_{\mathcal{D}}[f_{\mathcal{D}}(\mathbf{x})]$$

$$\bar{f}(\mathbf{x}) \approx \frac{1}{K} \sum_{k=1}^K f_{\mathcal{D}^{(k)}}(\mathbf{x})$$

K training sets (of size N) sampled from
 $P(\mathbf{x}, y): \mathcal{D}^{(1)}, \mathcal{D}^{(2)}, \dots, \mathcal{D}^{(K)}$

Using the average hypothesis

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) + \bar{f}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \end{aligned}$$

Using the average hypothesis

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) + \bar{f}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) \right)^2 + \left(\bar{f}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \end{aligned}$$

Using the average hypothesis

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] &= \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) + \bar{f}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) \right)^2 + \left(\bar{f}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right]\end{aligned}$$

Bias and variance

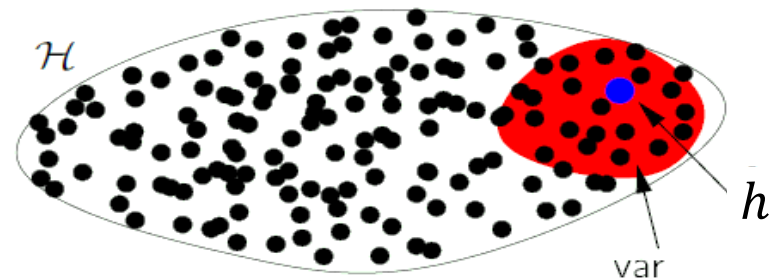
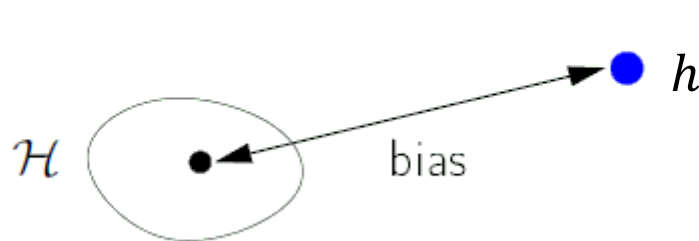
$$\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] = \underbrace{\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 \right]}_{\text{var}(\mathbf{x})} + \underbrace{(\bar{f}(\mathbf{x}) - h(\mathbf{x}))^2}_{\text{bias}(\mathbf{x})}$$

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] \right] &= \mathbb{E}_{\mathbf{x}} [\text{var}(\mathbf{x}) + \text{bias}(\mathbf{x})] \\ &= \text{var} + \text{bias} \end{aligned}$$

Example: sin target

$$\text{var} = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}) \right)^2 \right] \right]$$

$$\text{bias} = \mathbb{E}_{\mathbf{x}} [\bar{f}(\mathbf{x}) - h(\mathbf{x})]$$

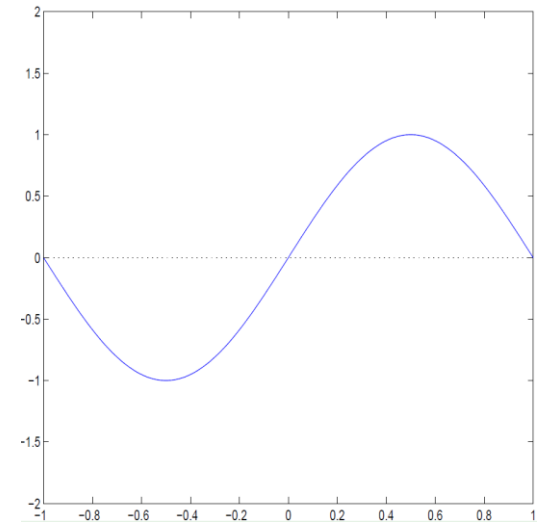


More complex $\mathcal{H} \Rightarrow$ lower bias but higher variance

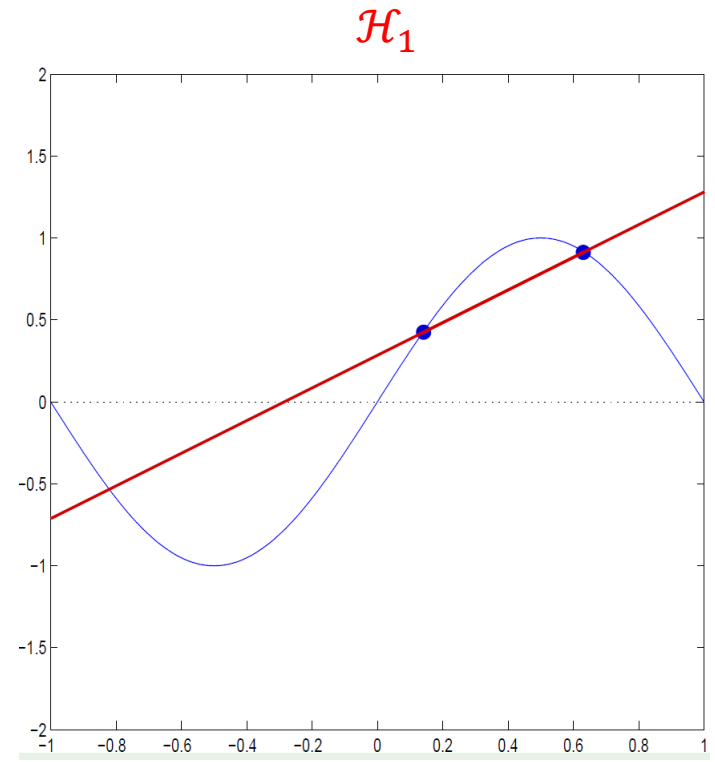
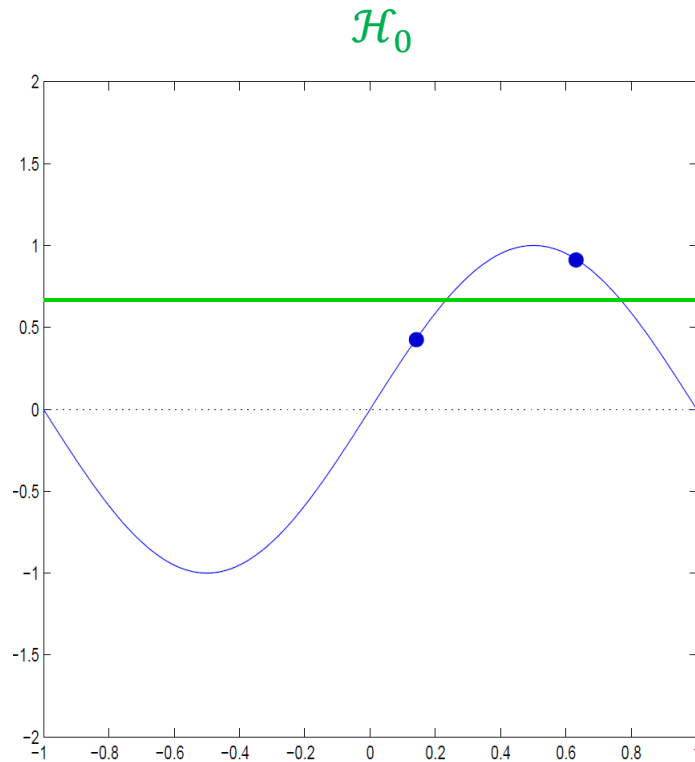
[Y.S. Abou Mostafa, 2012]

Example: sin target

- ▶ Only two training example $N = 2$
- ▶ Two models used for learning:
 - ▶ $\mathcal{H}_0: f(x) = b$
 - ▶ $\mathcal{H}_1: f(x) = ax + b$
- ▶ Which is better \mathcal{H}_0 or \mathcal{H}_1 ?

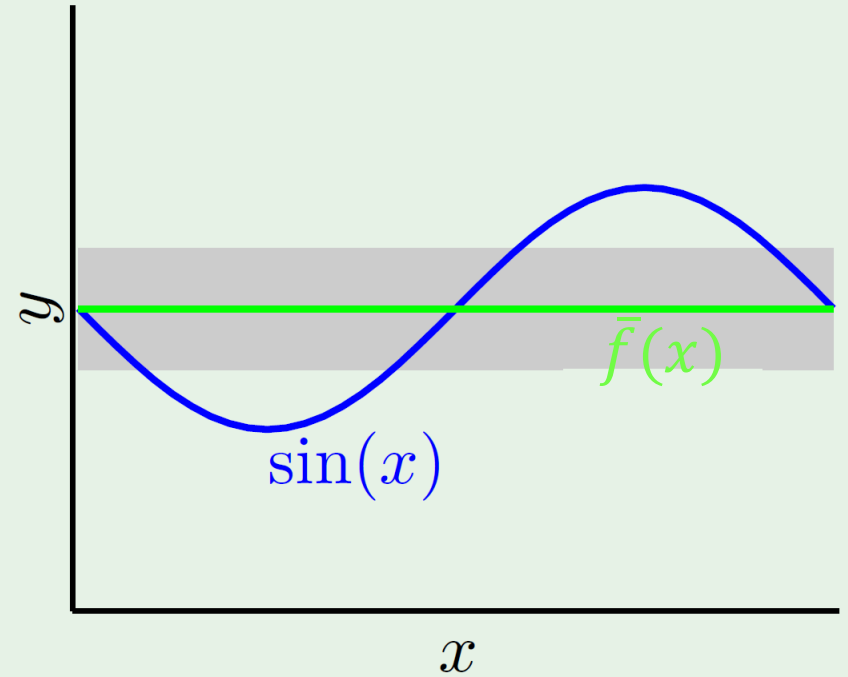
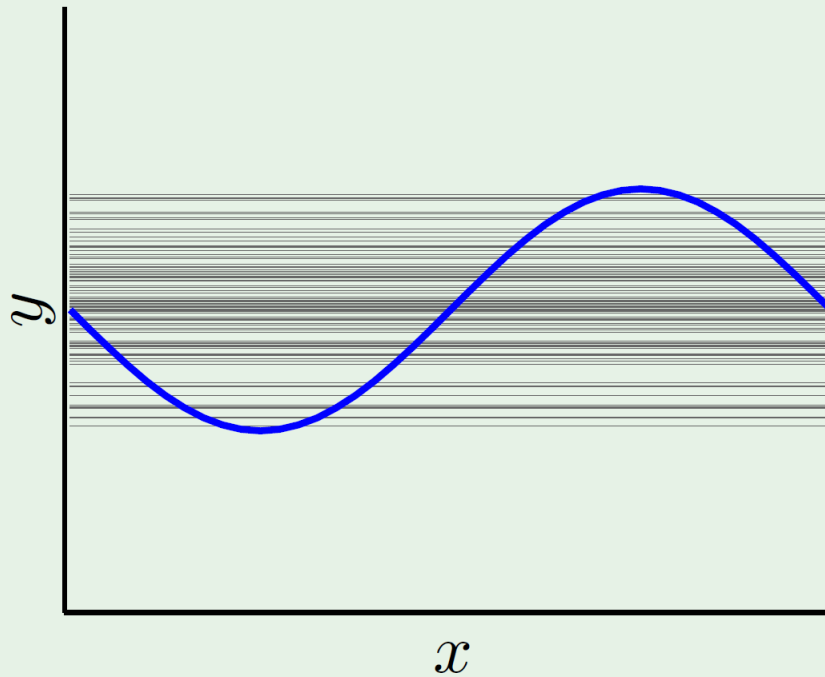


Learning from a training set



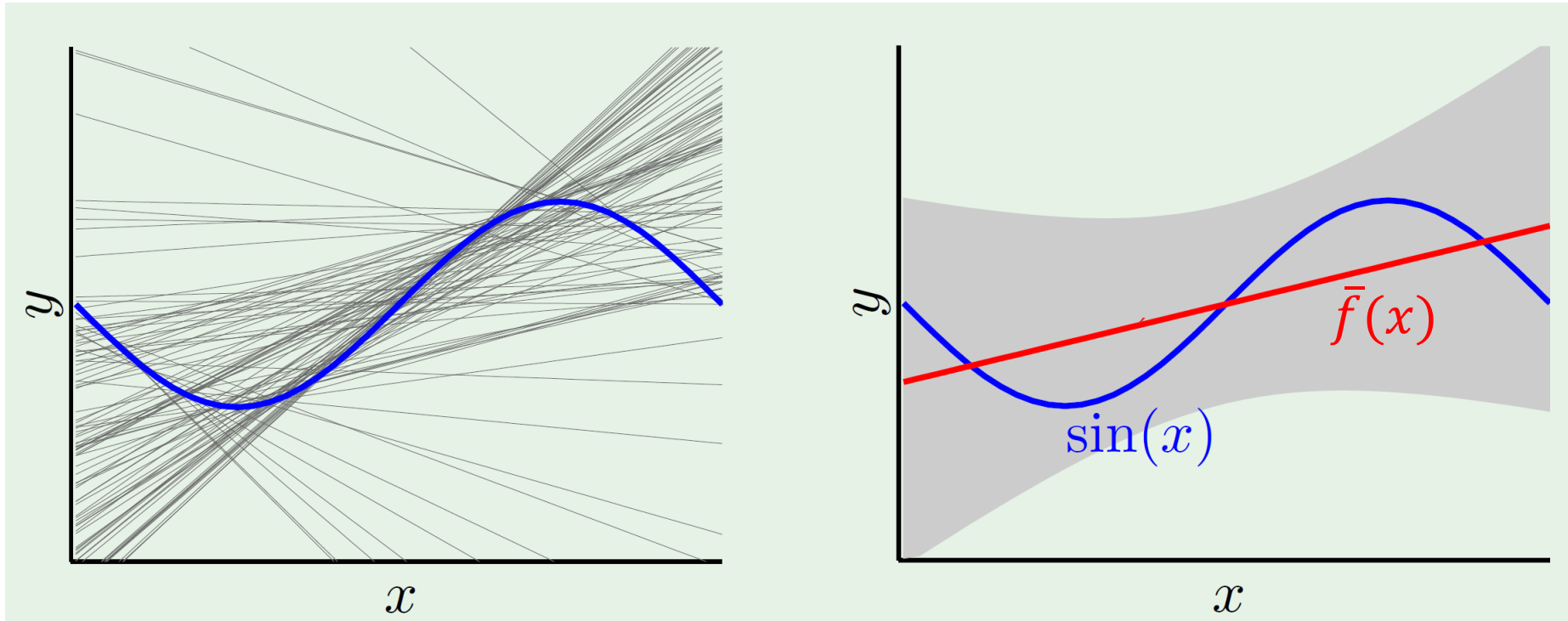
[Y.S. Abou Mostafa, 2012]

Variance \mathcal{H}_0



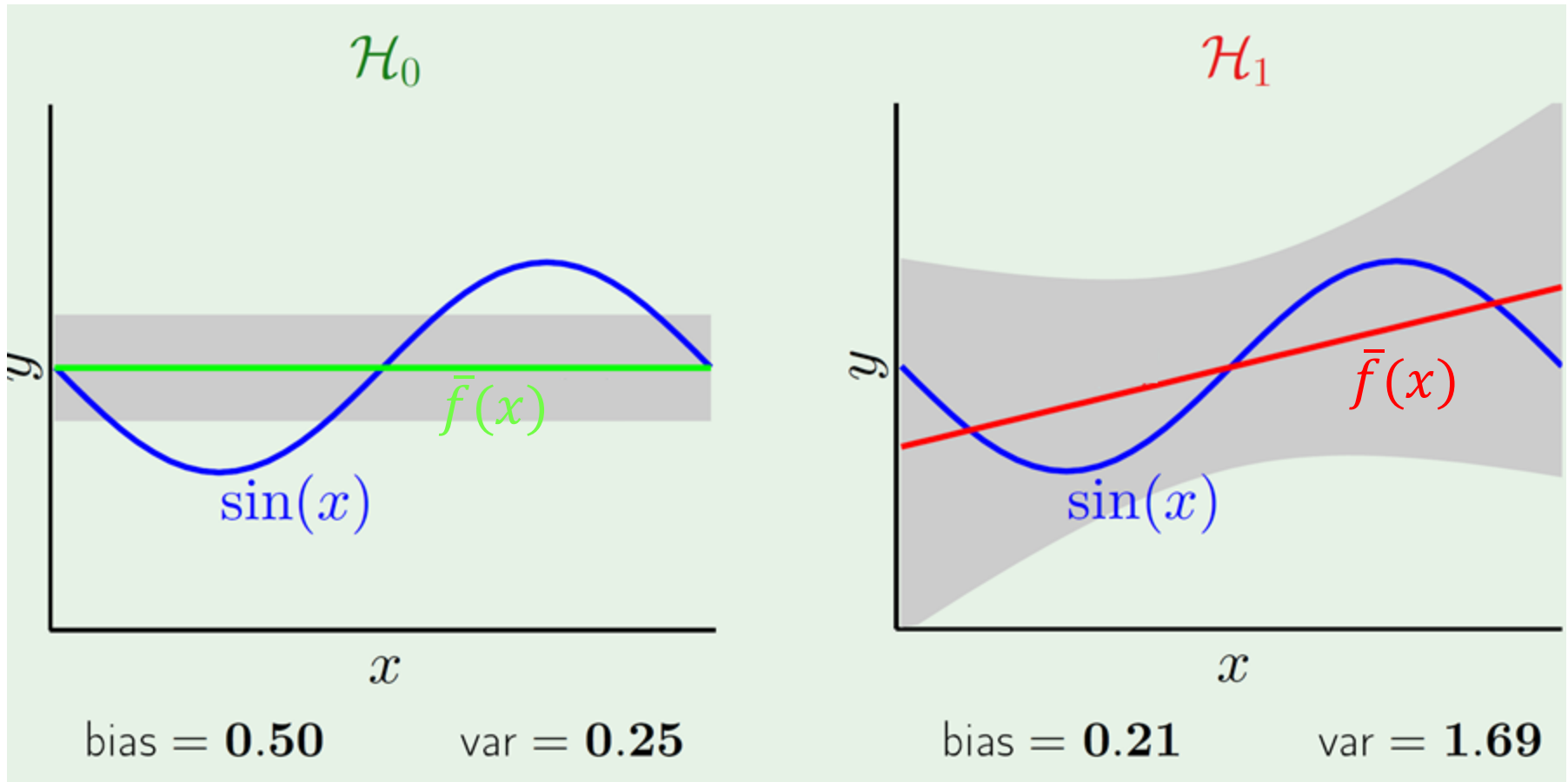
[Y.S. Abou Mostafa, et. al]

Variance \mathcal{H}_1



[Y.S. Abou Mostafa, et. al]

Which is better?



[Y.S. Abou Mostafa, 2012]

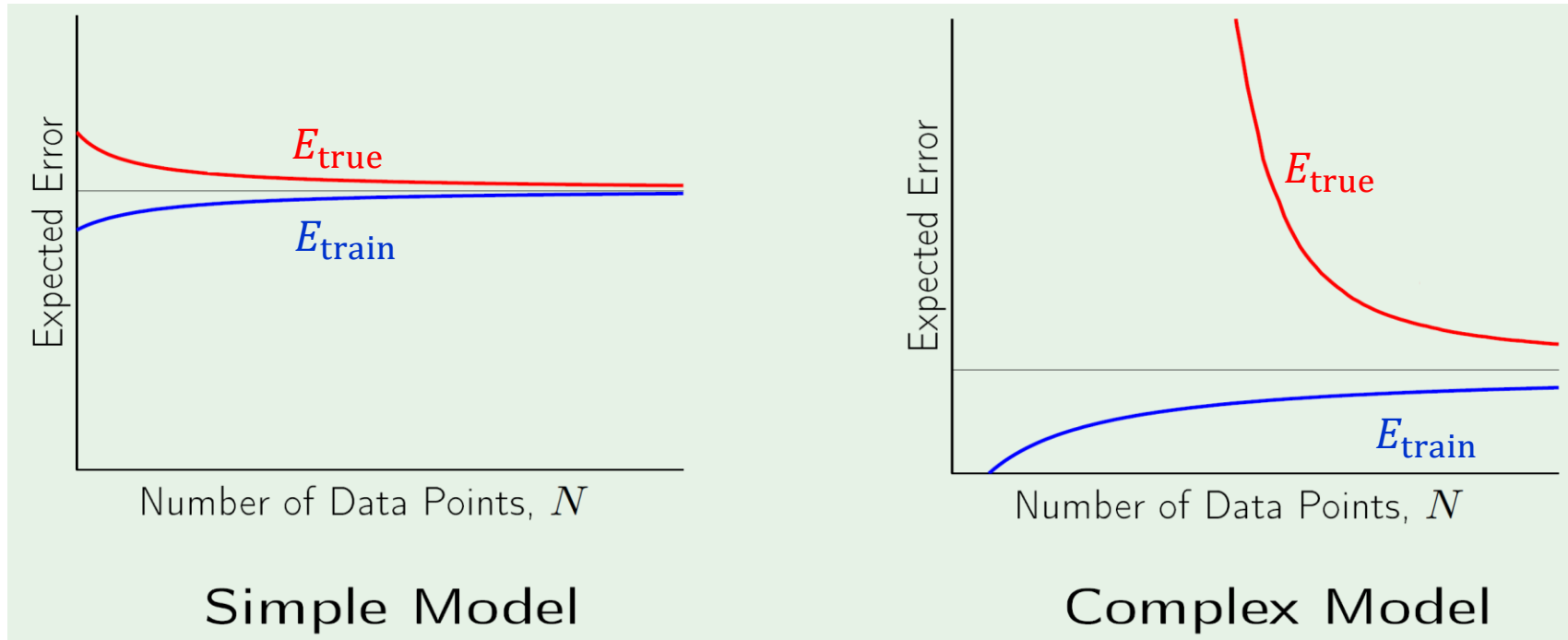
Lesson

Match the **model complexity**
to the **data sources**
not to the complexity of the **target function**.



Expected training and true error curves

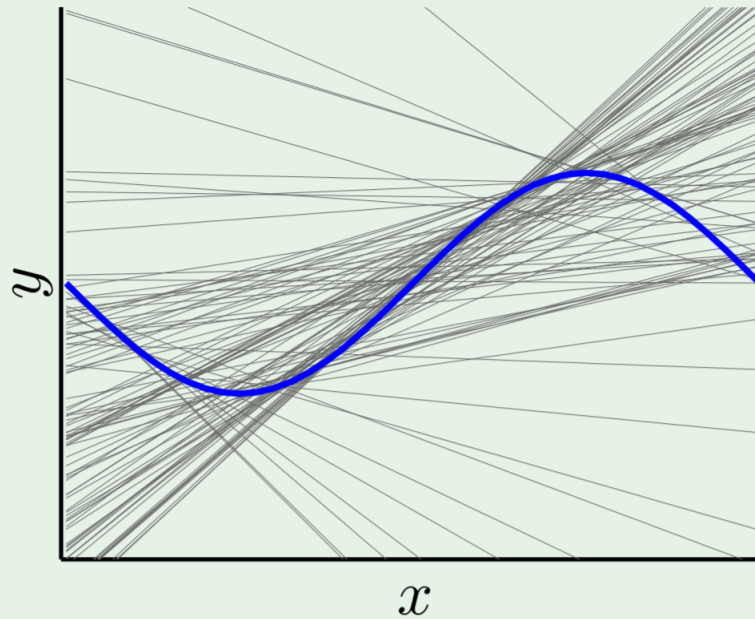
- Errors vary with the number of training samples



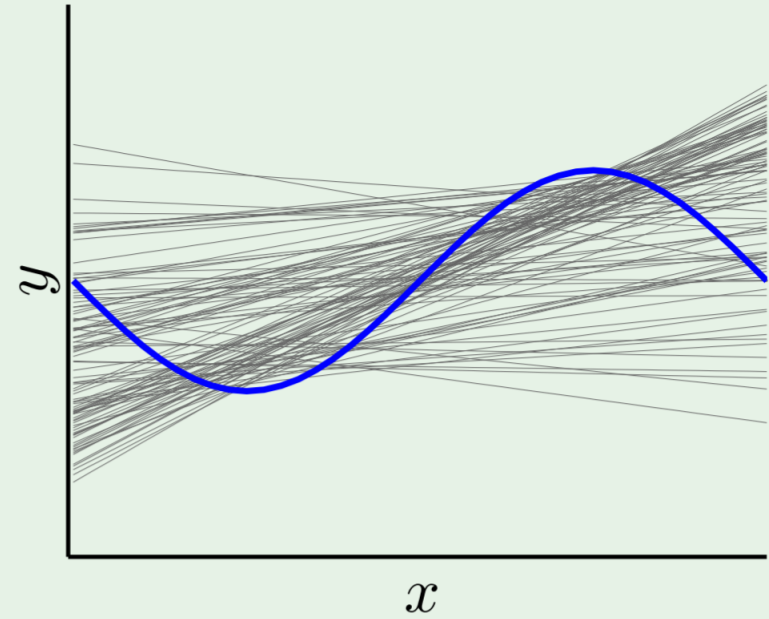
expected true error: $\mathbb{E}_{\mathcal{D}}[E_{\text{true}}(f_{\mathcal{D}}(\mathbf{x}))]$
expected training error: $\mathbb{E}_{\mathcal{D}}[E_{\text{train}}(f_{\mathcal{D}}(\mathbf{x}))]$

[Y.S. Abou Mostafa, 2012]

Regularization



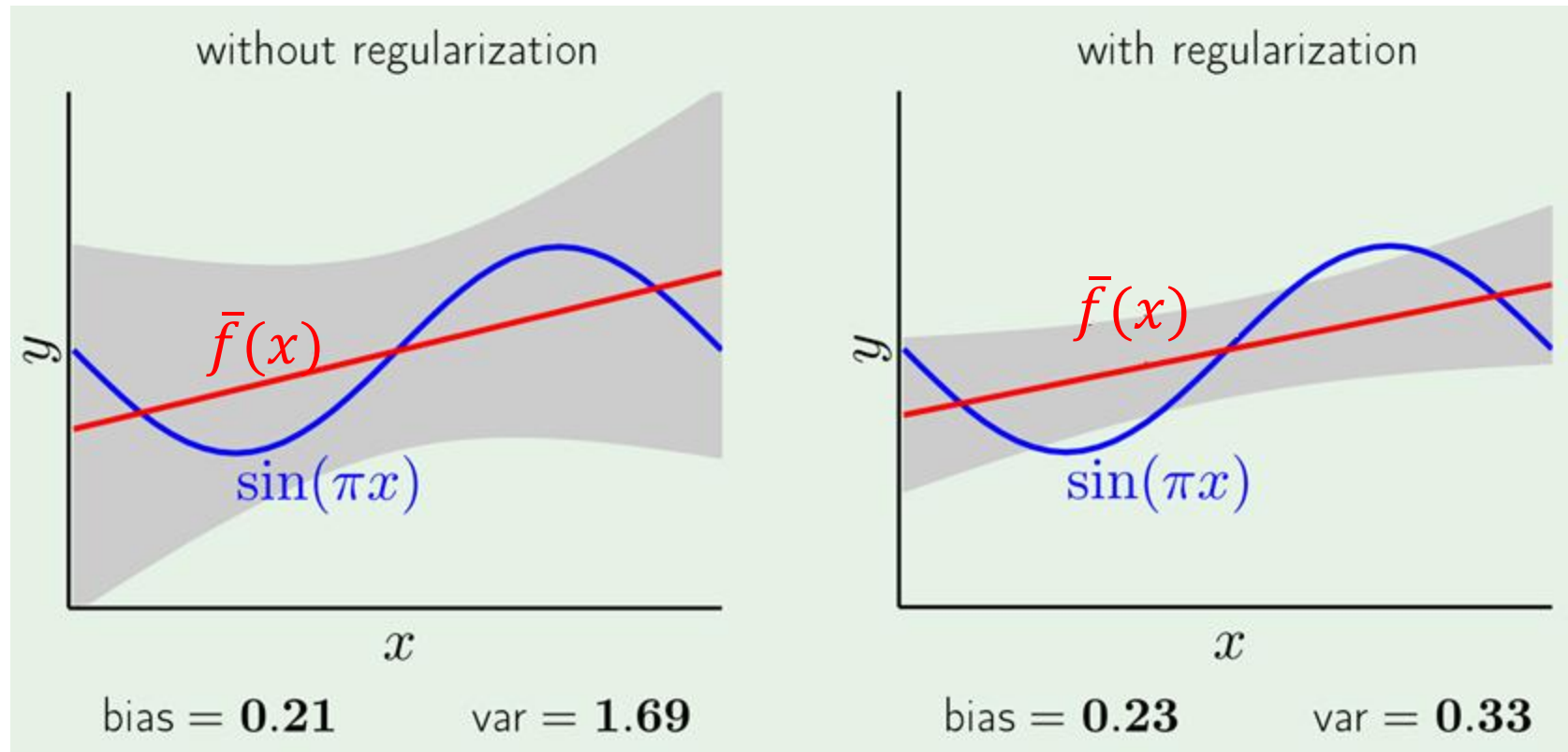
without regularization



with regularization

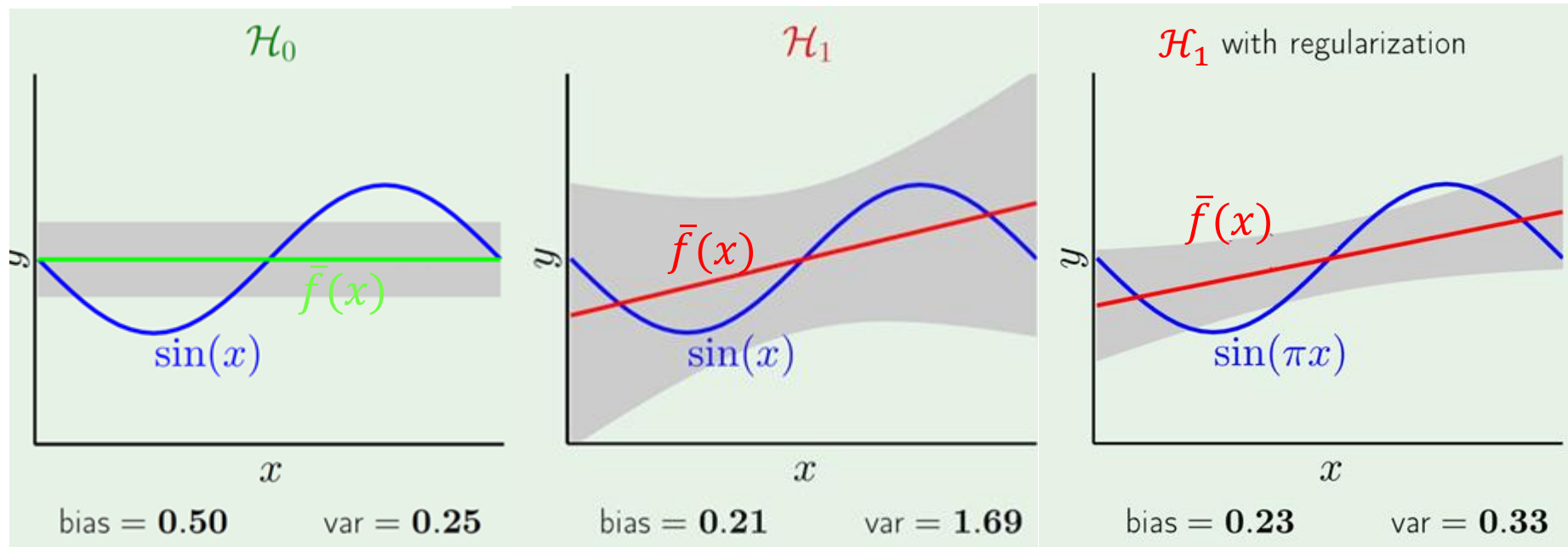
[Y.S. Abou Mostafa, 2012]

Regularization: bias and variance



[Y.S. Abou Mostafa, 2012]

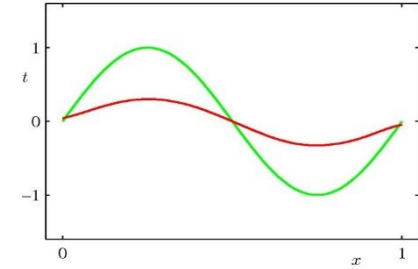
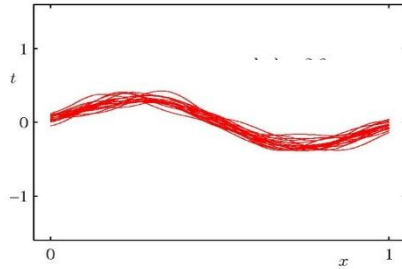
Resource Winner of \mathcal{H}_0 , \mathcal{H}_1 , and \mathcal{H}_1 with regularization



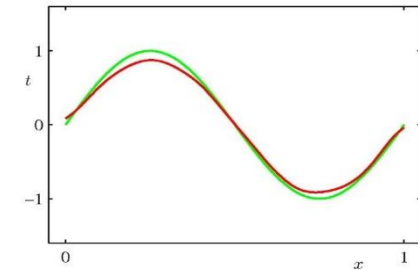
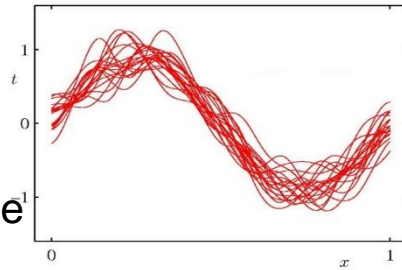
[Y.S. Abou Mostafa, 2012]

Regularization and bias/variance

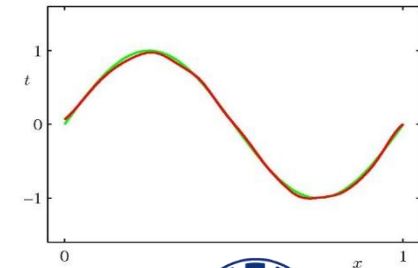
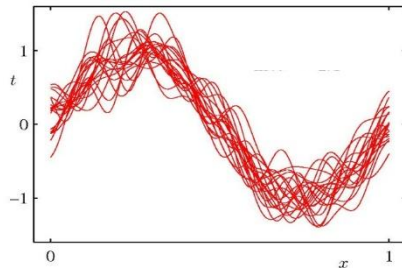
λ is large



λ is intermediate



λ is small



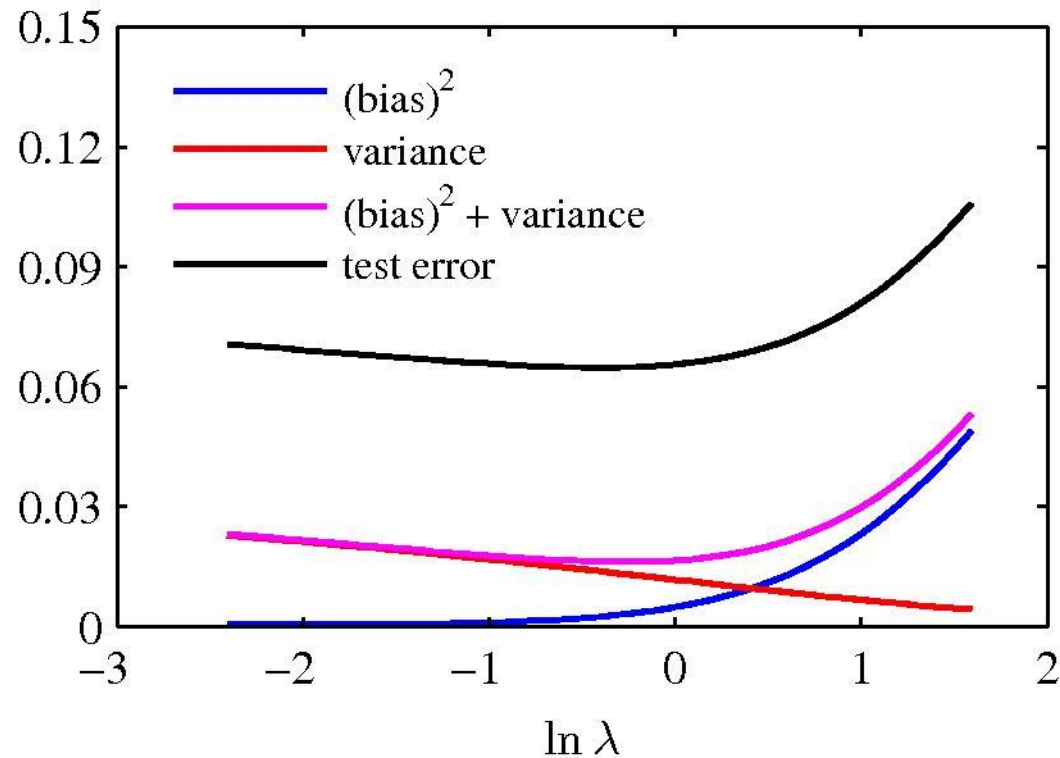
$L = 100$ data sets

$N = 25$

$m = 25$

[Bishop]

Learning curves of bias, variance, and noise



[Bishop]

Bias-variance decomposition: summary

- ▶ The noise term is unavoidable.
- ▶ The terms we are interested in are bias and variance.
- ▶ The approximation-generalization trade-off is seen in the bias-variance decomposition.

Resources

- ▶ C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 3.2.
- ▶ Yaser S. Abu-Mostafa, Malik Maghdon-Ismael, and Hsuan Tien Lin, “**Learning from Data**”, Chapter 2.3, 3.2, 3.4.
- ▶ Course CE-717, Dr. M.Soleyman

