



Generalization, Regularization Bias-variance trade-off

Iran University of Science and Technology

By: M. S. Tahaei Ph.D.

Fall 2024

Courtesy: slides are adopted partly from Dr. Soleymani, Sharif University

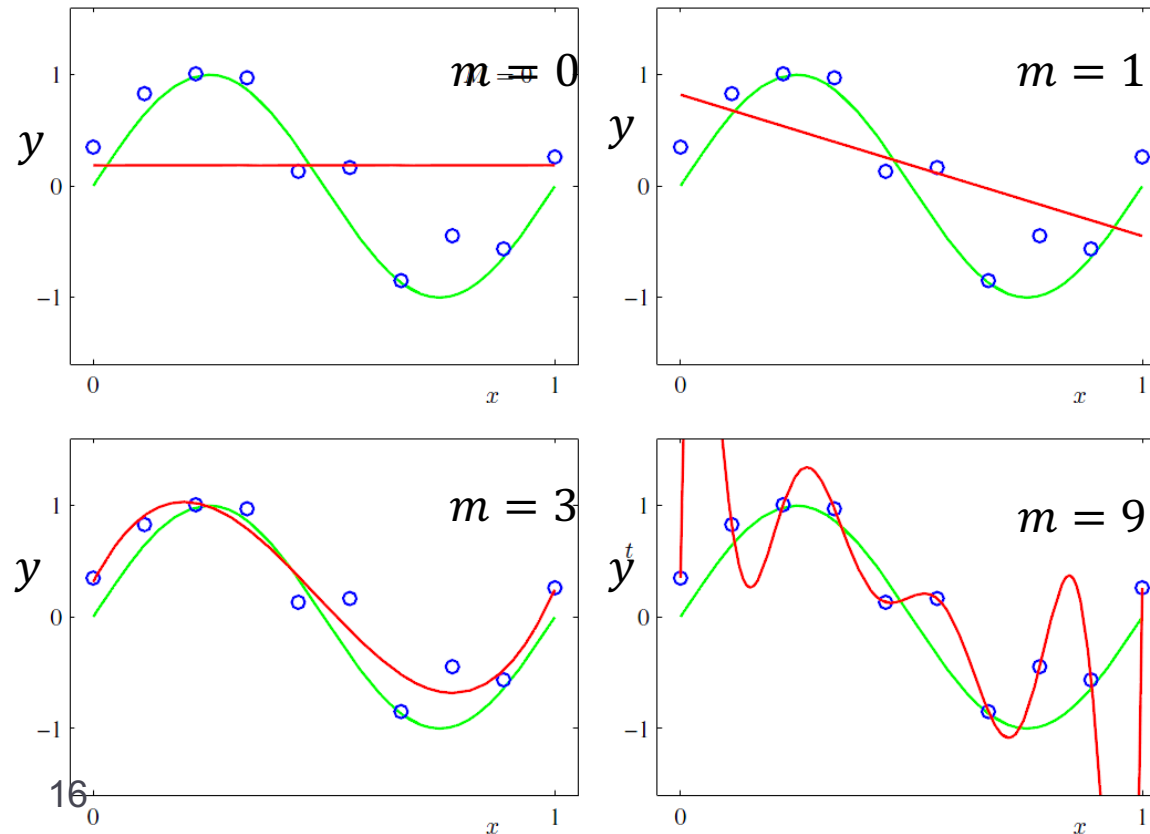
Outline

- Generalization
- Regularization
- Bias-variance

Model complexity

- ▶ Example:

- ▶ Polynomials with larger m are becoming increasingly tuned to the random noise on the target values.



Evaluation and model selection

- ▶ **Evaluation:**

- ▶ We need to measure how well the learned function can predicts the target for unseen examples

- ▶ **Model selection:**

- ▶ Most of the time we need to select among a set of models
 - ▶ Example: polynomials with different degree m
- ▶ and thus we need to evaluate these models first

Avoiding over-fitting

- ▶ Determine a suitable value for model complexity
 - ▶ **Simple hold-out method**
 - ▶ **Cross-validation**
- ▶ Regularization (Occam's Razor)
 - ▶ Explicit preference towards simple models
 - ▶ Penalize for the model complexity in the objective function

Simple hold-out: model selection

- ▶ Steps:

- ▶ Divide training data into training and validation set v_set
- ▶ Use only the training set to train a set of models
- ▶ Evaluate each learned model on the validation set

- ▶ $J_v(\mathbf{w}) = \frac{1}{|v_set|} \sum_{i \in v_set} \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$

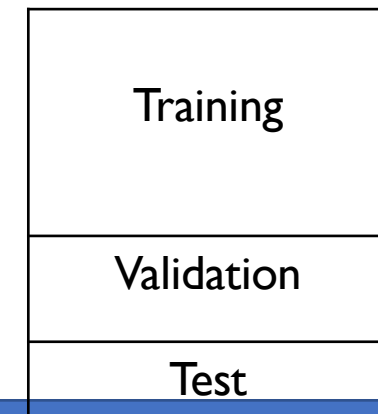
- ▶ Choose the best model based on the validation set error

- ▶ Usually, too wasteful of valuable training data

- ▶ Training data may be limited.
 - ▶ On the other hand, small validation set give a relatively noisy estimate of performance.

Simple hold out: training, validation, and test sets

- ▶ Simple hold-out chooses the model that minimizes error on validation set.
- ▶ $J_v(\mathbf{w})$ is likely to be an optimistic estimate of generalization error.
 - ▶ extra parameter (e.g., degree of polynomial) is fit to this set.
- ▶ Estimate generalization error for the test set
 - ▶ performance of the selected model is finally evaluated on the test set



Cross-Validation (CV): Evaluation

- ▶ k -fold cross-validation steps:
 - ▶ Shuffle the dataset and randomly partition training data into k groups of approximately equal size
 - ▶ for $i = 1$ to k
 - ▶ Choose the i -th group as the held-out validation group
 - ▶ Train the model on all but the i -th group of data
 - ▶ Evaluate the model on the held-out group
 - ▶ Performance scores of the model from k runs are **averaged**.
 - ▶ The average error rate can be considered as an estimation of the true performance.

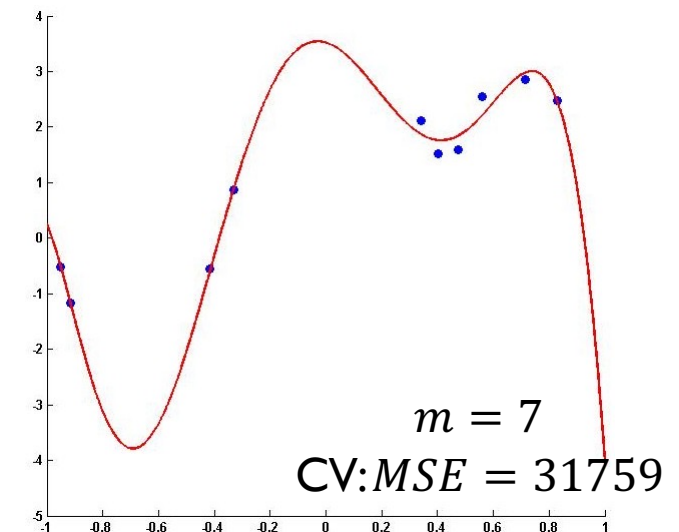
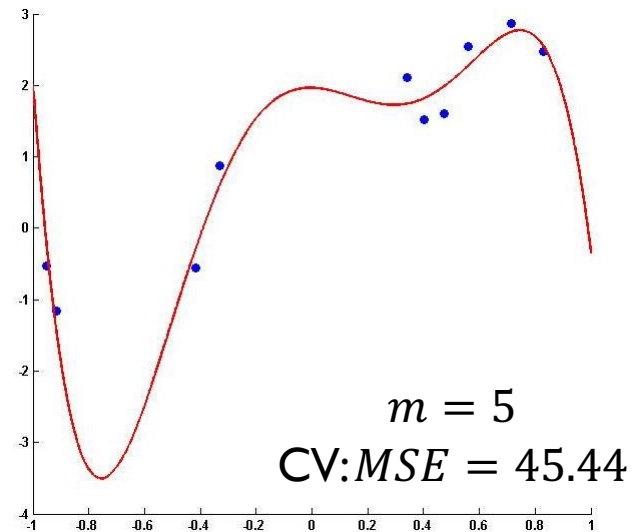
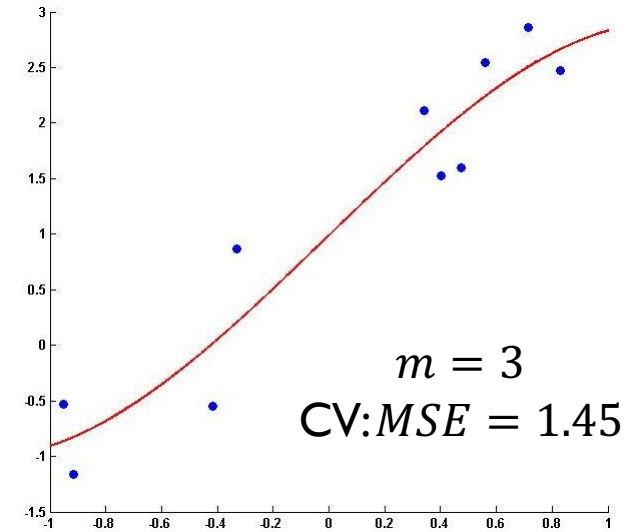
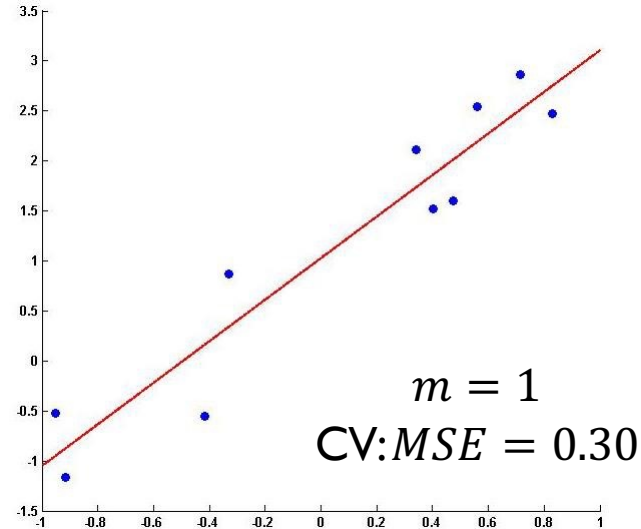


Cross-Validation (CV): Model Selection

- ▶ For each model we first find the average error find by CV.
- ▶ The model with the best average performance is selected.

Cross-validation: polynomial regression example

5-fold CV
100 runs
average



Leave-One-Out Cross Validation (LOOCV)

- ▶ When data is particularly scarce, cross-validation with $k = N$
 - ▶ Leave-one-out treats each training sample in turn as a test example and all other samples as the training set.
- ▶ Use for small datasets
 - ▶ When training data is valuable
 - ▶ LOOCV can be time expensive as N training steps are required.

Regularization

- ▶ Adding a penalty term in the cost function to discourage the coefficients from reaching large values.
- ▶ Ridge regression (weight decay):

$$J(\mathbf{w}) = \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(x^{(i)}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w}$$

$$\mathbf{w} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

Polynomial order

- ▶ Polynomials with larger m are becoming increasingly tuned to the random noise on the target values.
- ▶ magnitude of the coefficients typically gets larger by increasing m .

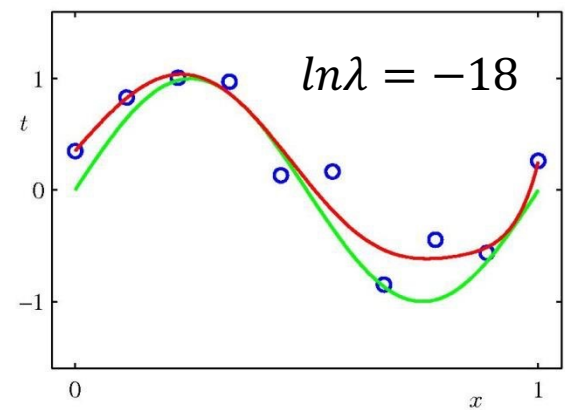
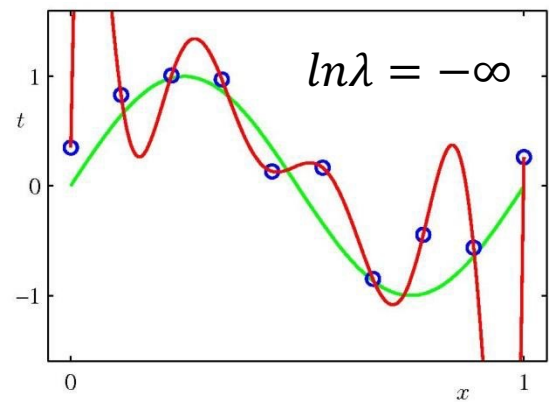
	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

[Bishop]

Regularization parameter

	0	$m = 9$	1
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0	0.35	0.35	0.13
w_1	232.37	4.74	-0.05
w_2	-5321.83	-0.77	-0.06
w_3	48568.31	-31.97	-0.05
w_4	-231639.30	-3.89	-0.03
w_5	640042.26	55.28	-0.02
w_6	-1061800.52	41.32	-0.01
w_7	1042400.18	-45.95	-0.00
w_8	-557682.99	-91.53	0.00
w_9	125201.43	72.68	0.01

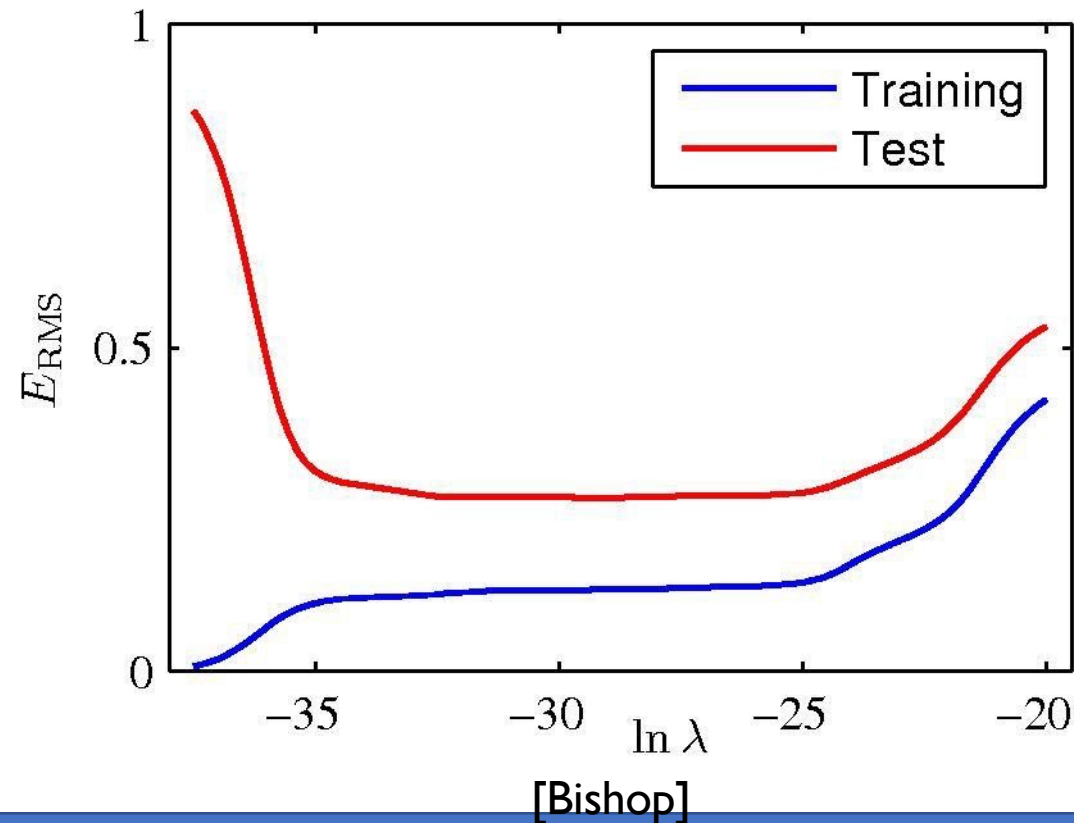
[Bishop]



Regularization parameter

- ▶ Generalization

- ▶ λ now controls the effective complexity of the model and hence determines the degree of over-fitting



The approximation-generalization trade-off

- ▶ Small true error shows good approximation of f out of sample
- ▶ More complex $\mathcal{H} \Rightarrow$ better chance of approximating f
- ▶ Less complex $\mathcal{H} \Rightarrow$ better chance of generalization out of f

Complexity of Hypothesis Space: Example

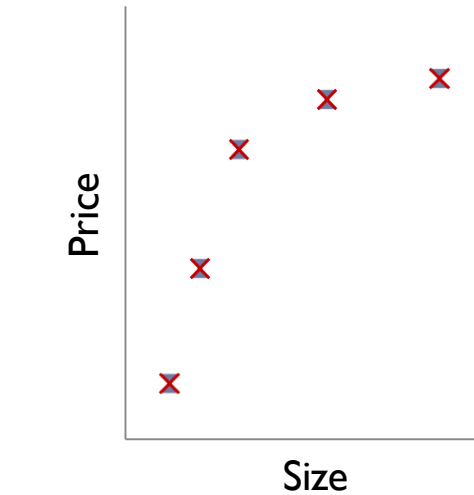


$$w_0 + w_1x$$

Less complex \mathcal{H}



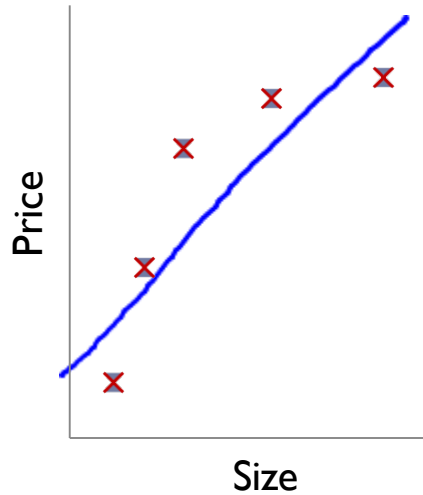
$$w_0 + w_1x + w_2x^2$$



$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

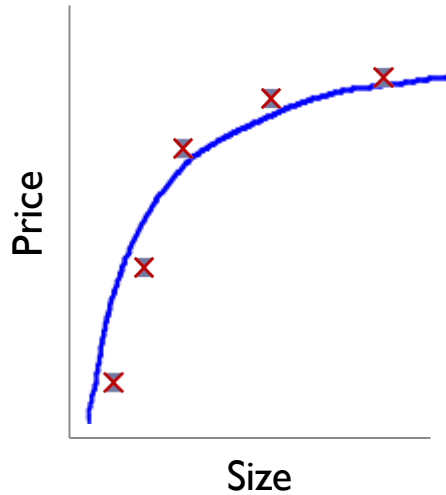
More complex \mathcal{H}

Complexity of Hypothesis Space: Example

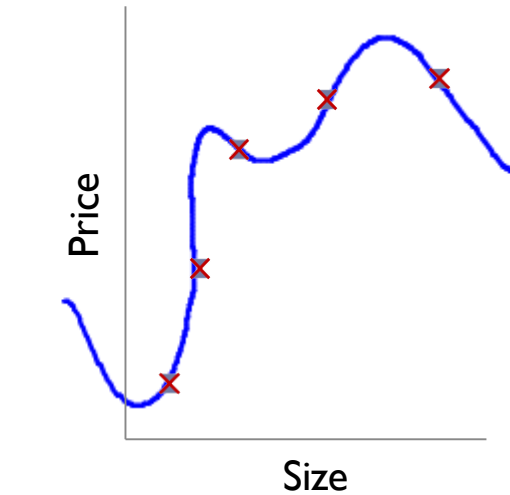


$$w_0 + w_1x$$

Less complex \mathcal{H}



$$w_0 + w_1x + w_2x^2$$

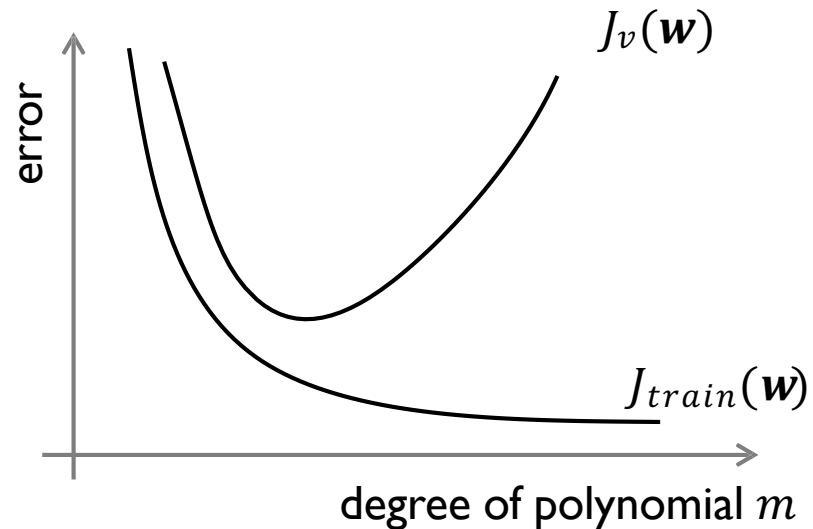


$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

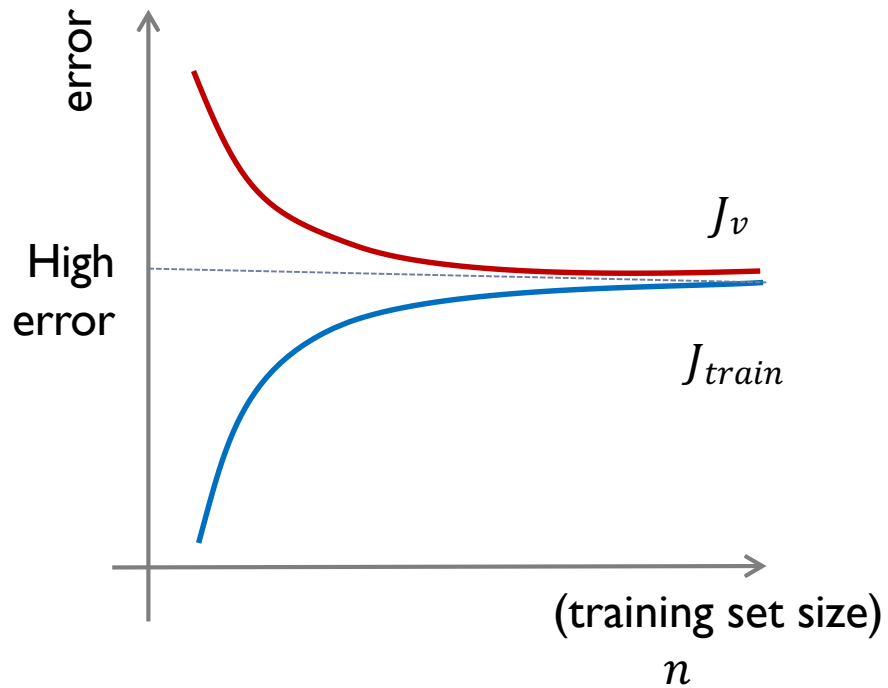
More complex \mathcal{H}

Complexity of Hypothesis Space

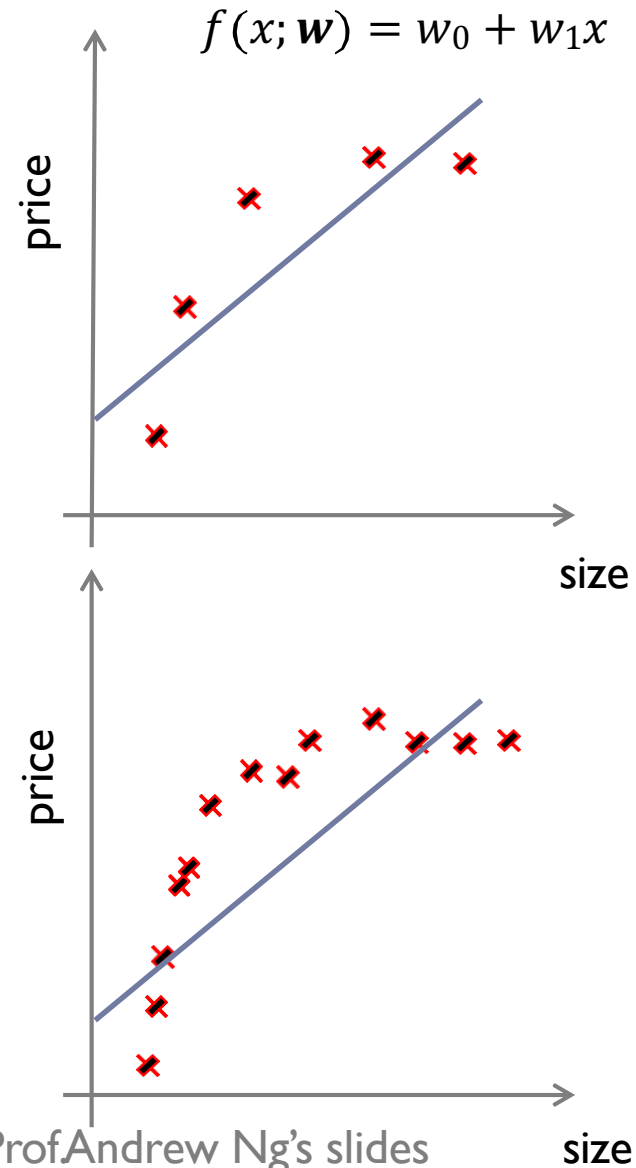
- ▶ Less complex \mathcal{H} :
 - ▶ $J_{train}(\mathbf{w}) \approx J_v(\mathbf{w})$ and $J_{train}(\mathbf{w})$ is very high
- ▶ More complex \mathcal{H} :
 - ▶ $J_{train}(\mathbf{w}) \ll J_v(\mathbf{w})$ and $J_{train}(\mathbf{w})$ is low



Less complex \mathcal{H}

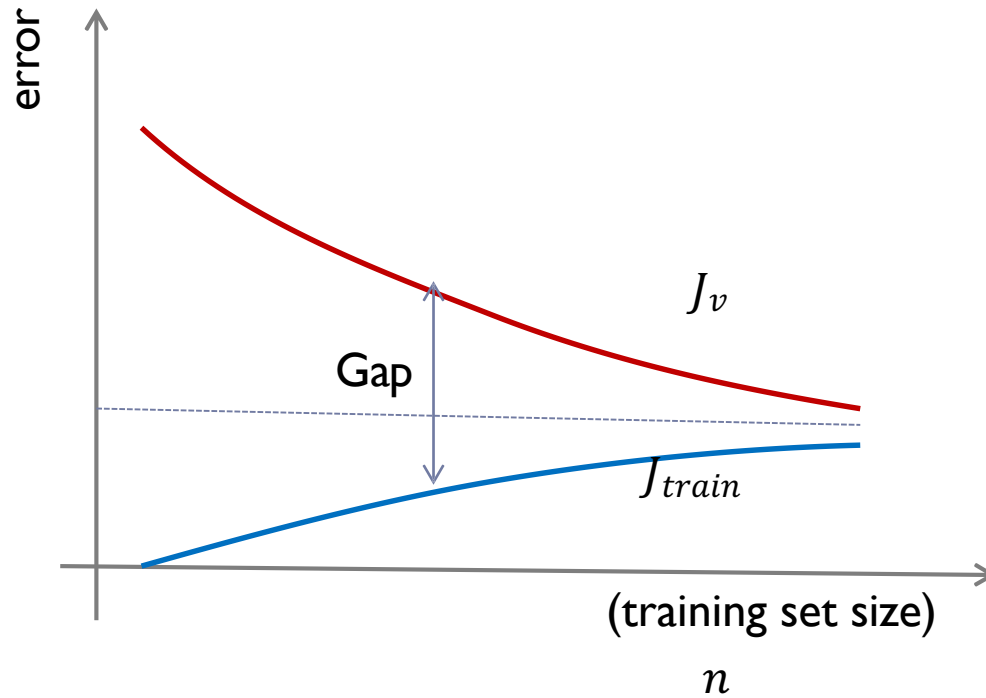


If model is very simple, getting more training data will not (by itself) help much.

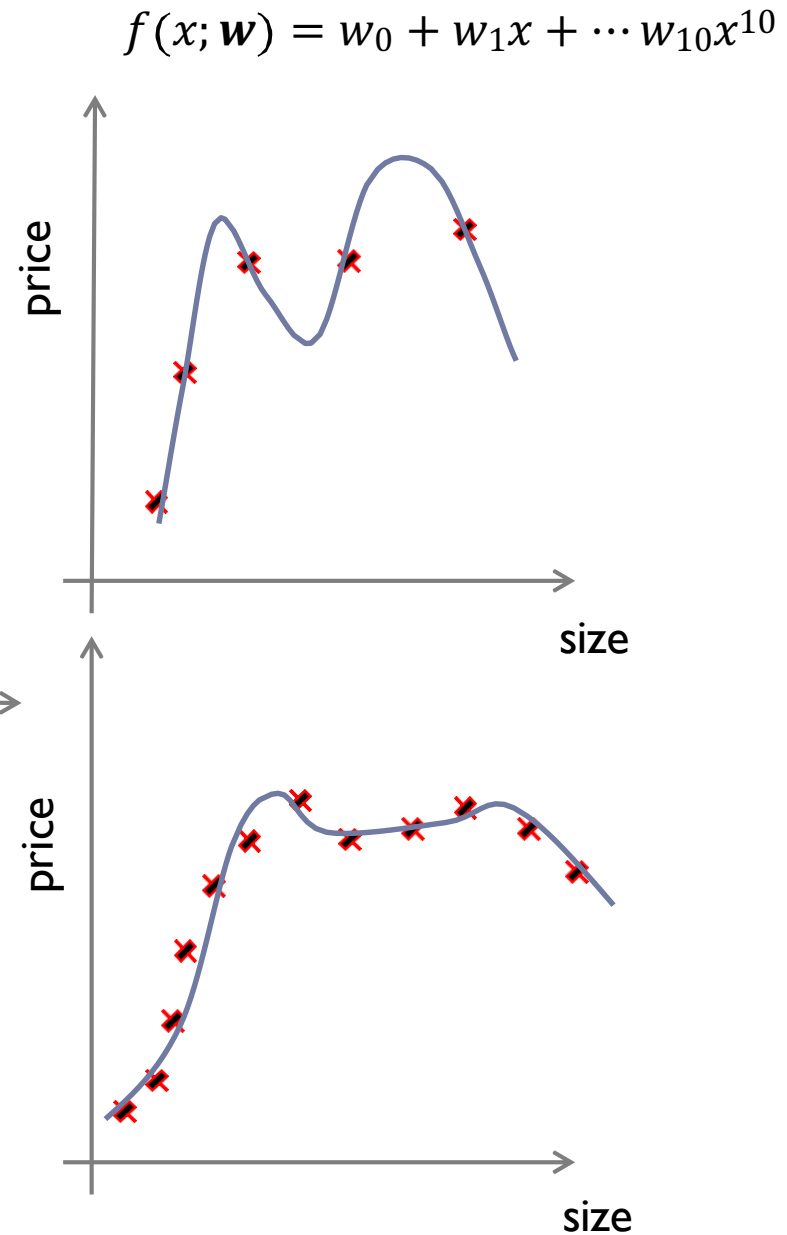


This slide has been adapted from: Prof Andrew Ng's slides

More complex \mathcal{H}



For more complex models, getting more training data is usually helps.

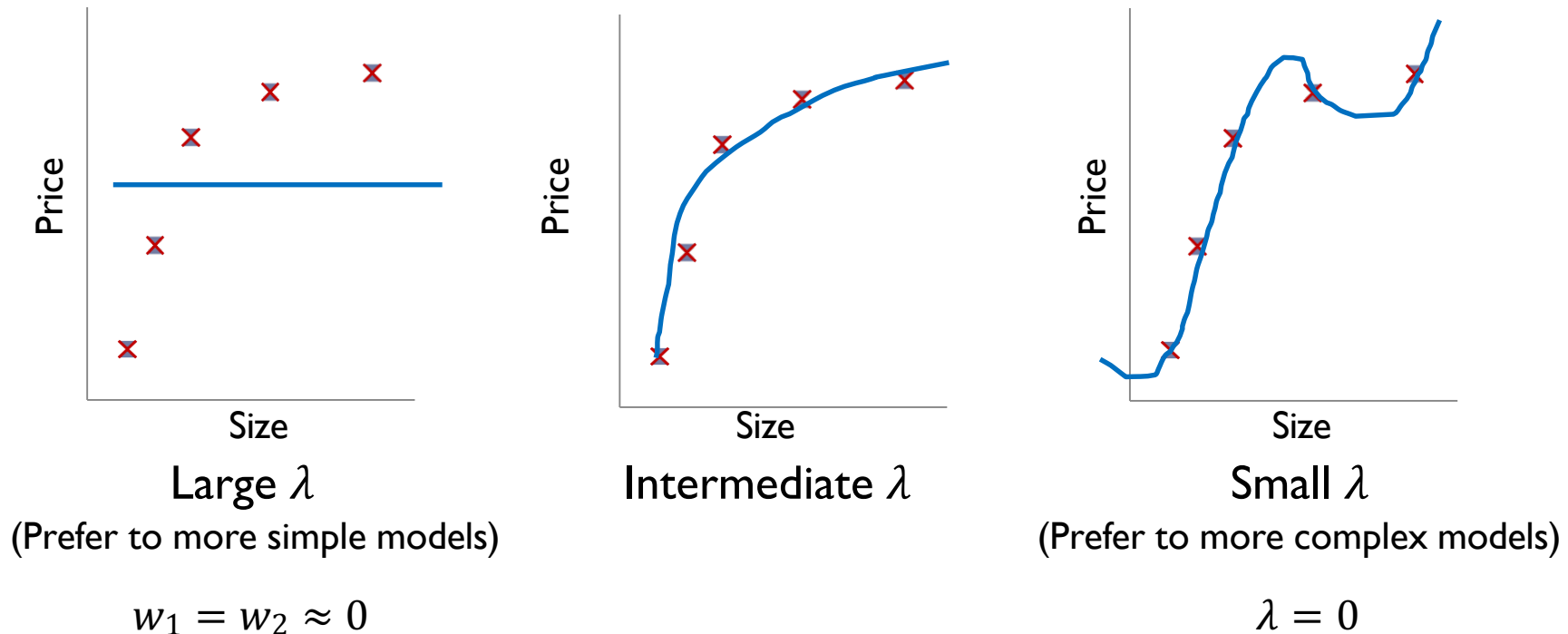


This slide has been adapted from: Prof Andrew Ng's slides

Regularization: Example

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

$$J(\mathbf{w}) = \frac{1}{n} \left(\sum_{i=1}^n \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w} \right)$$



This example has been adapted from: Prof. Andrew Ng's slides

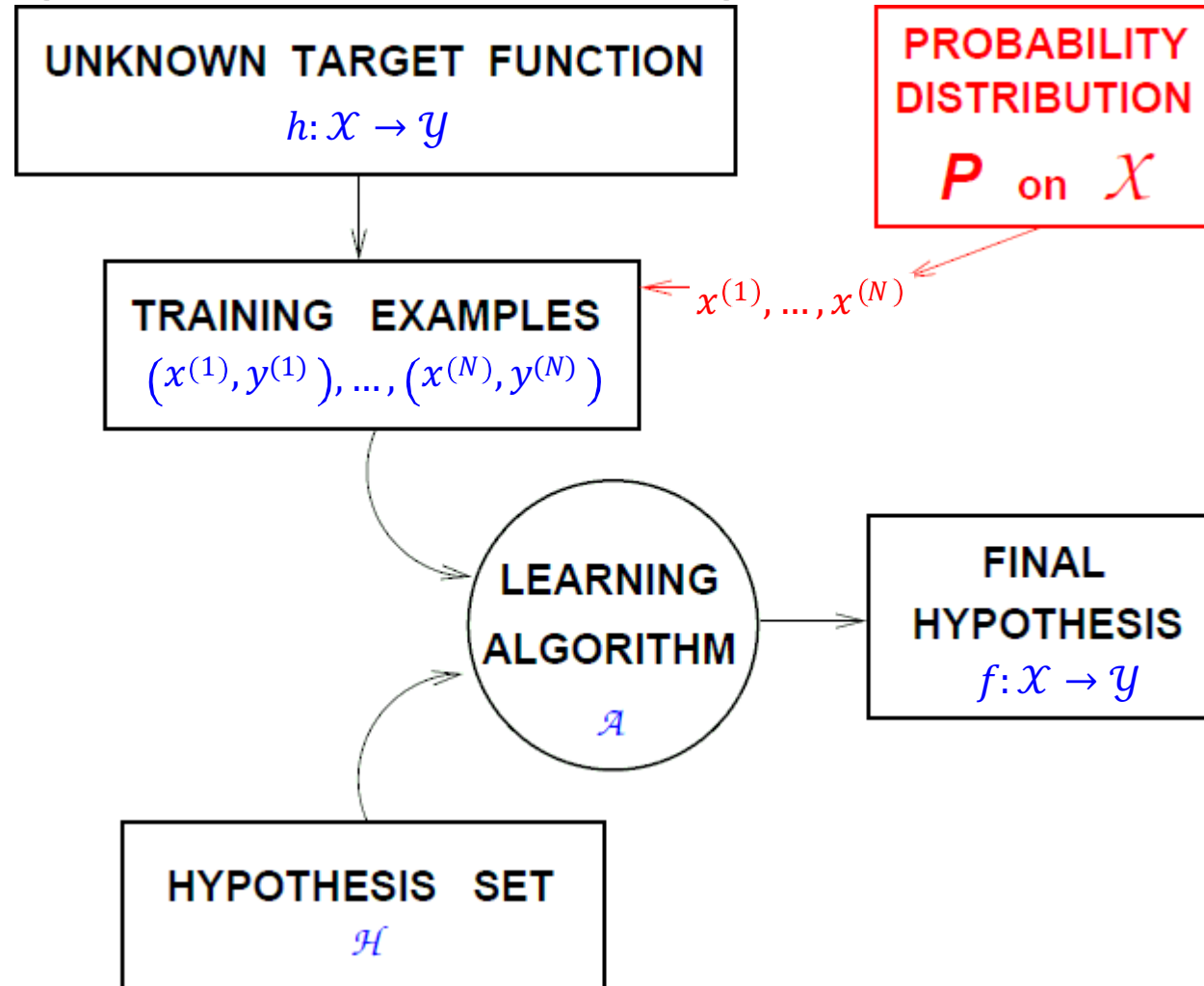
Model complexity: Bias-variance trade-off

- ▶ 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

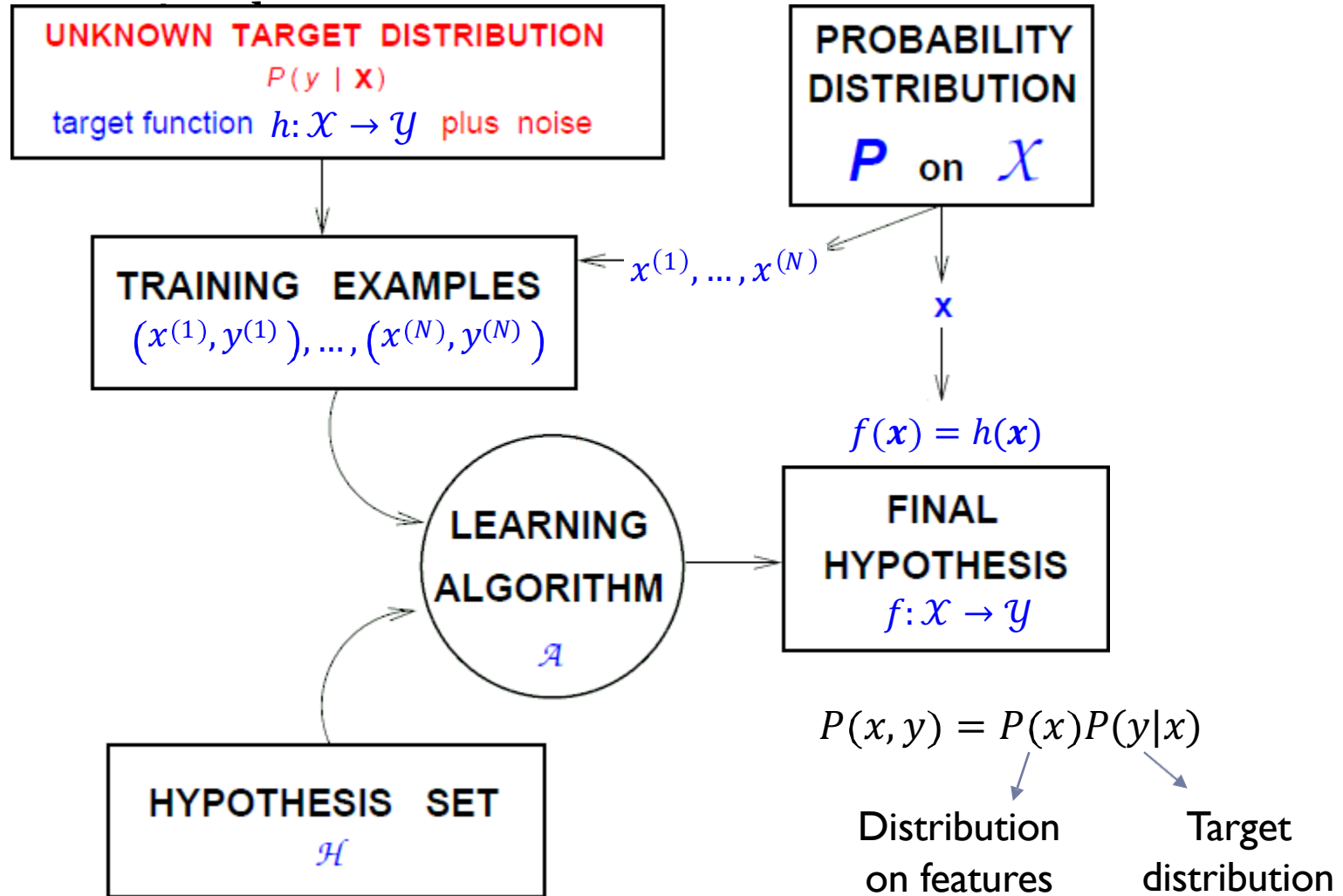
- ▶ Best unrestricted regression function
- ▶ Noise
- ▶ Bias and variance

The learning diagram: deterministic target



The learning diagram including noisy target

► Type



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

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

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

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

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

Error decomposition

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

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

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

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

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

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

0

Error decomposition

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

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

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

- Noise shows the irreducible minimum value of the loss function

Expectation of true error

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

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

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

The average hypothesis

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

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

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

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

Bias and variance

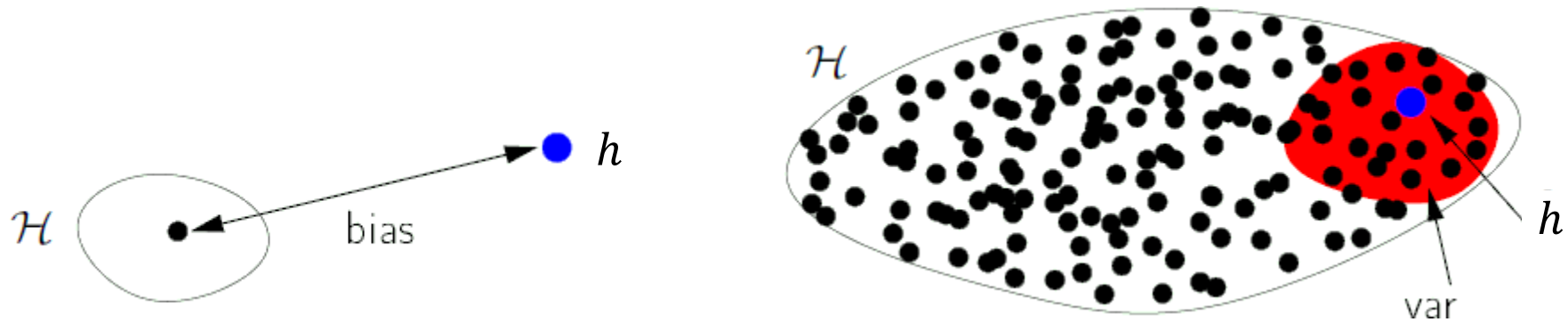
$$\mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}) \right)^2 \right] = \underbrace{\mathbb{E}_{\mathcal{D}} \left[\left(f_{\mathcal{D}}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]}_{\text{var}(\mathbf{x})} + \underbrace{\left(f(\mathbf{x}) - h(\mathbf{x}) \right)^2}_{\text{bias}(\mathbf{x})}$$

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

Bias-variance trade-off

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

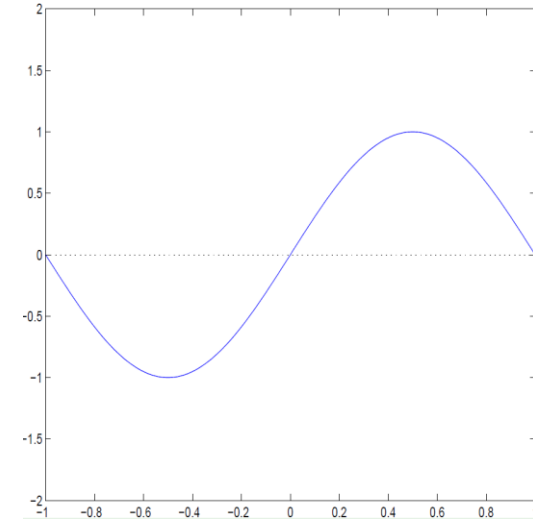
$$\text{bias} = \mathbb{E}_{\mathbf{x}} [f(\mathbf{x}) - h(\mathbf{x})]$$



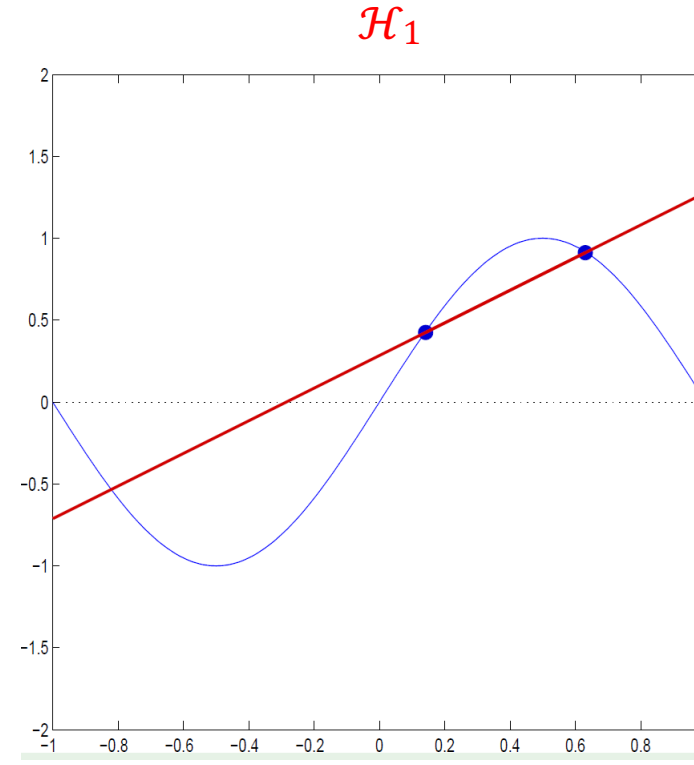
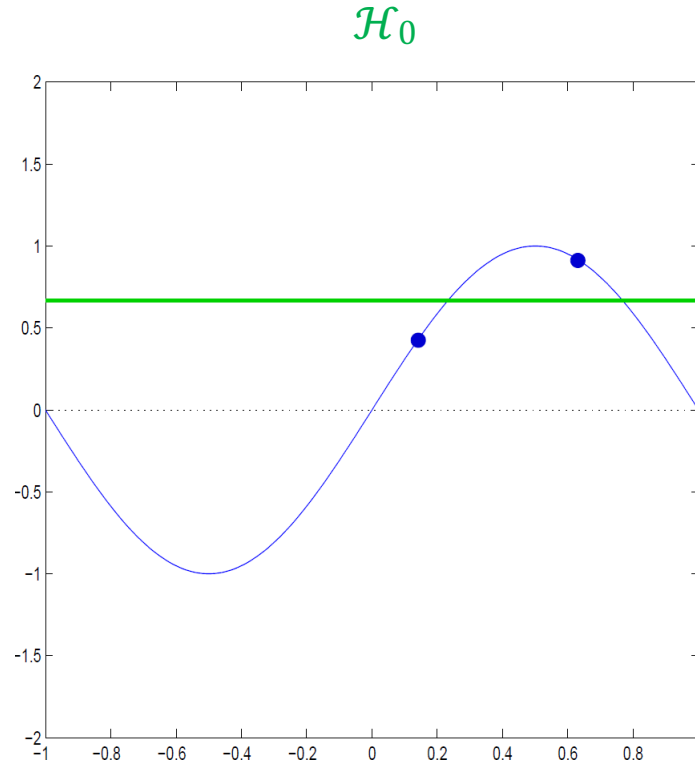
More complex $\mathcal{H} \Rightarrow$ lower bias but higher variance

Example: sin target

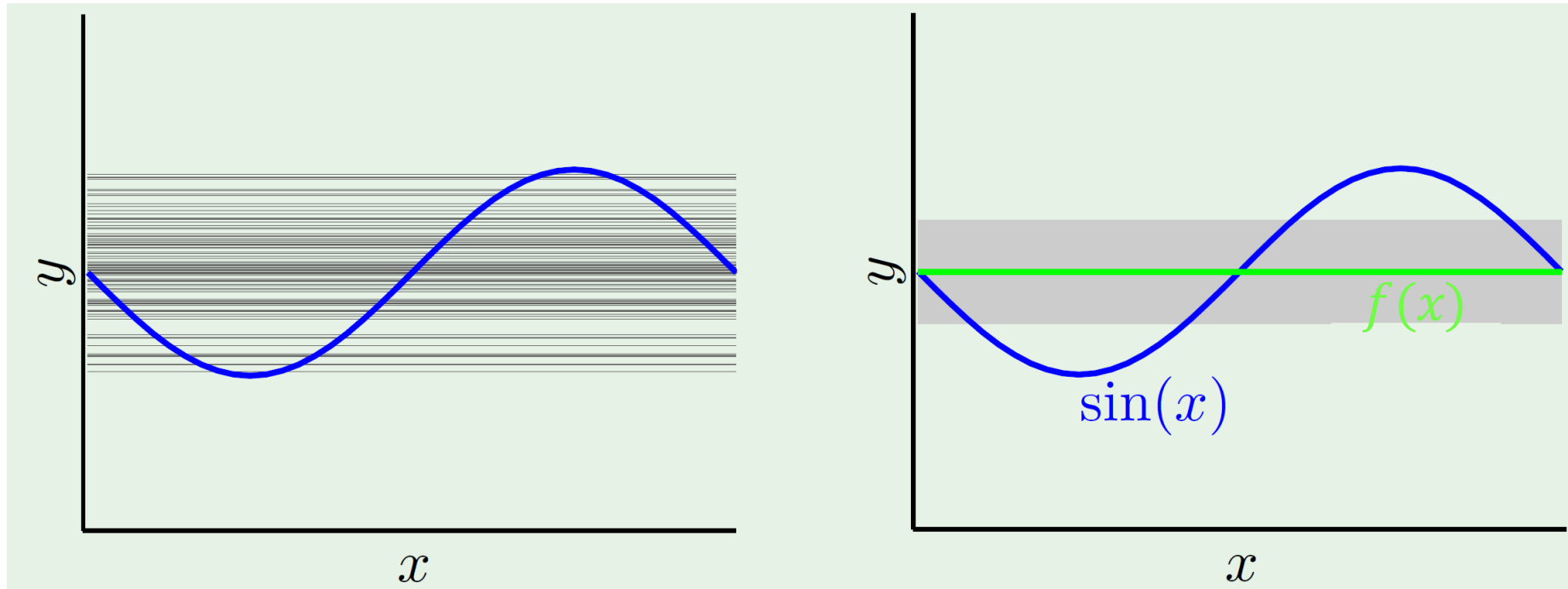
- ▶ Only two training example $N = 2$
- ▶ Two models used for learning:
 - ▶ $\mathcal{H}_0: f(x) = b$
 - ▶ $\mathcal{H}_1: f(x) = ax + b$
- ▶ Which is better \mathcal{H}_0 or \mathcal{H}_1 ?



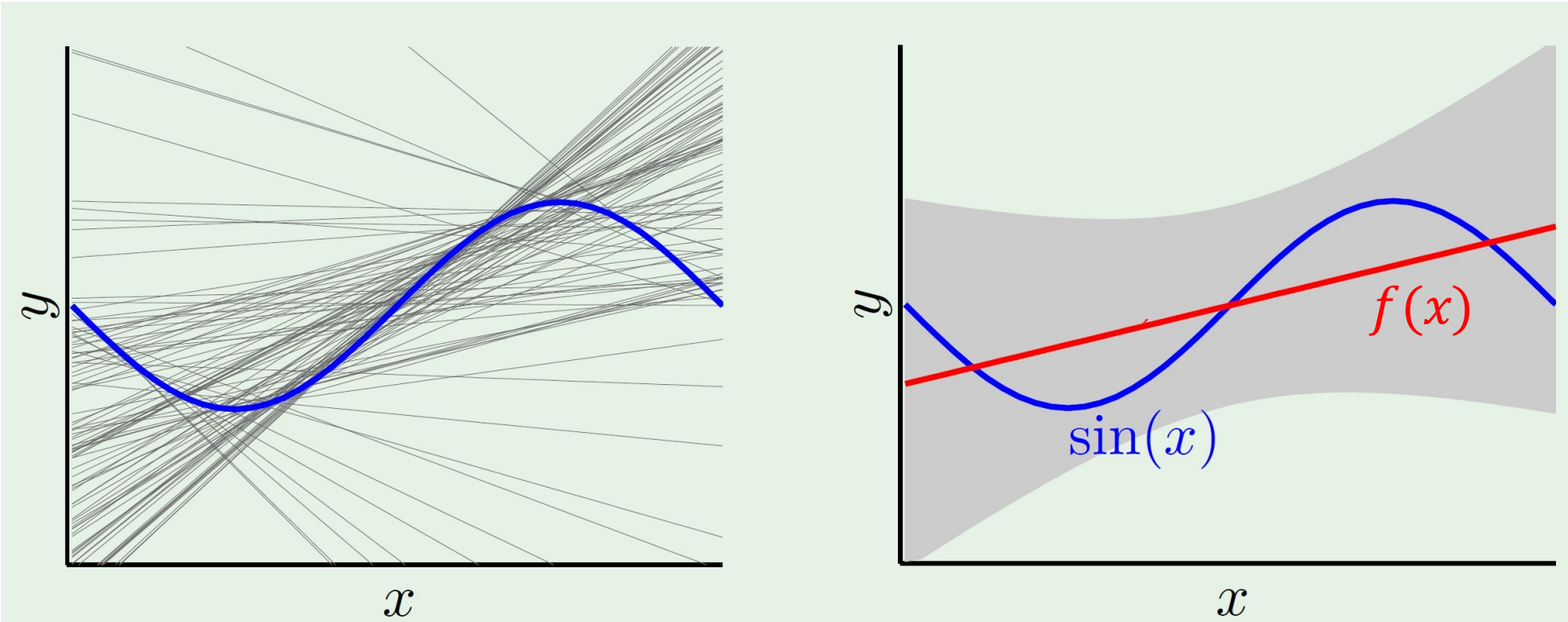
Learning from a training set



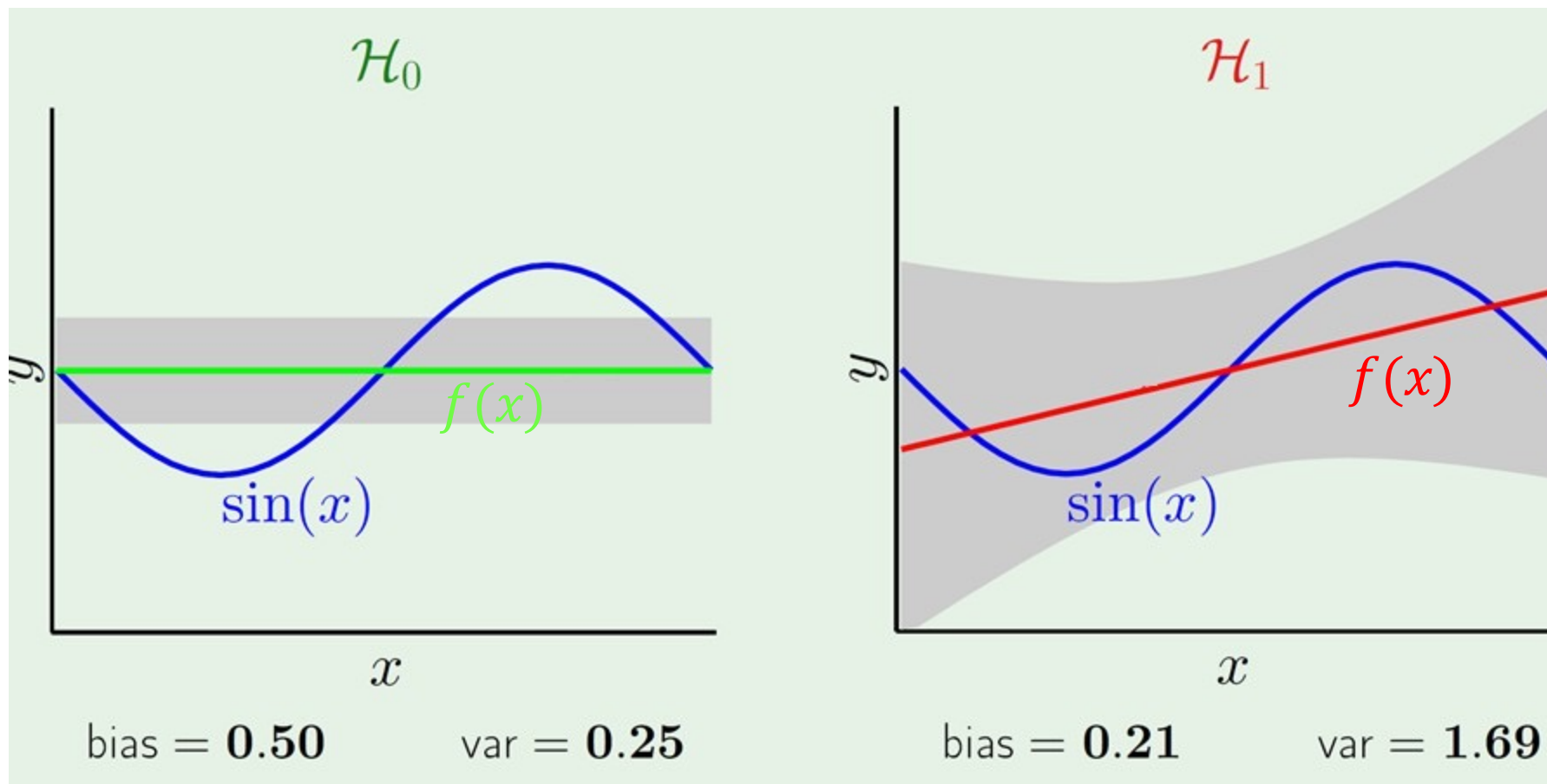
Variance \mathcal{H}_0



Variance \mathcal{H}_1



Which is better?



Resource

- 1 C. M. Bishop, *Pattern Recognition and Machine Learning*.
- 2 Y. S. Abu-Mostafa, "Machine learning." California Institute of Technology, 2012.
- 3 R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*. 2001.