

# Generalization

CE417: Introduction to Artificial Intelligence  
Sharif University of Technology  
Fall 2023

Soleymani

# Topics

- Beyond linear models
- Evaluation & model selection
- Regularization

# Recall: Linear regression (squared loss)

- Linear regression functions

$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$  are the parameters we need to set.

$$g : \mathbb{R} \rightarrow \mathbb{R} \quad g(x; \mathbf{w}) = w_0 + w_1 x$$

$$g : \mathbb{R}^d \rightarrow \mathbb{R} \quad g(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

- $J(\mathbf{w})$ : Sum of squares error

$$J(\mathbf{w}) = \sum_{i=1}^n \left( y^{(i)} - g(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$

- Weight update rule for  $g(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$ :

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \eta \sum_{i=1}^n \left( y^{(i)} - \mathbf{w}^{tT} \mathbf{x}^{(i)} \right) \mathbf{x}^{(i)}$$

# Beyond linear regression

- How to extend the linear regression to non-linear functions?
  - Transform the data using basis functions
  - Learn a linear regression on the new feature vectors (obtained by basis functions)

# Generalized linear

- Linear combination of fixed non-linear function of the input vector

$$g(\mathbf{x}; \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + \dots + w_m \phi_m(\mathbf{x})$$

$\{\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})\}$ : set of basis functions (or features)

$$\phi_i(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$$

# Basis functions: examples

- Linear

If  $m = d$ ,  $\phi_i(\mathbf{x}) = x_i$ ,  $i = 1, \dots, d$ , then

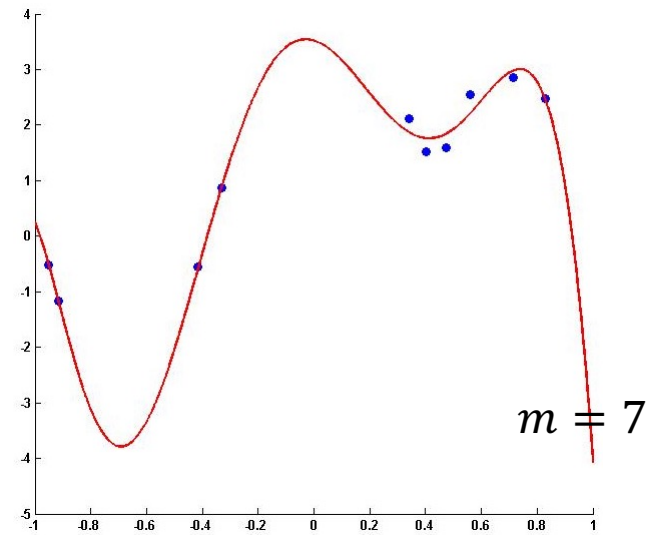
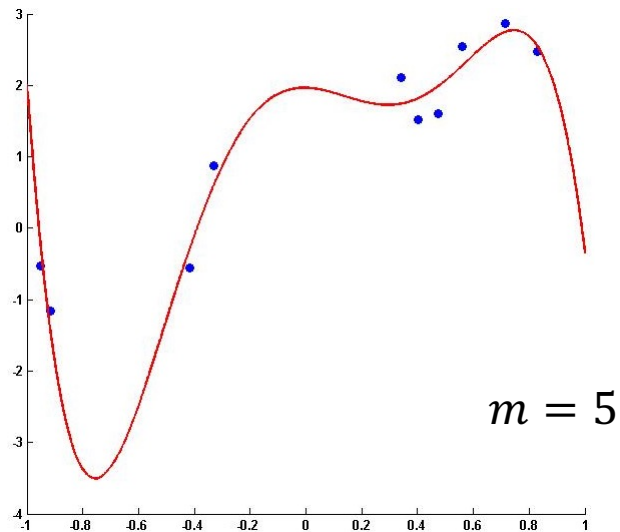
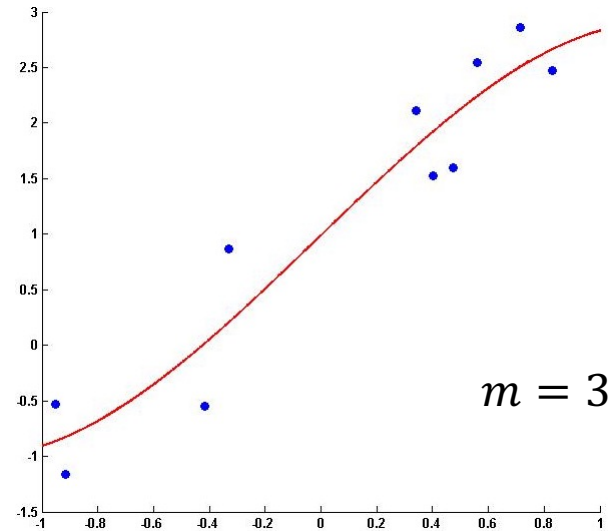
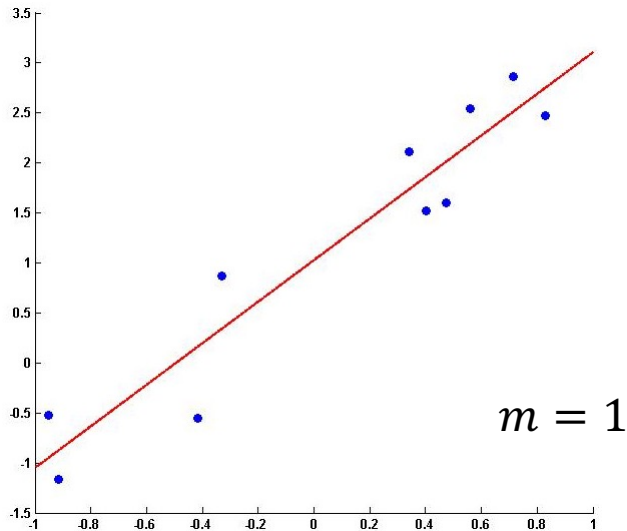
$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + \dots + w_dx_d$$

- Polynomial (univariate)

If  $\phi_i(x) = x^i$ ,  $i = 1, \dots, m$ , then

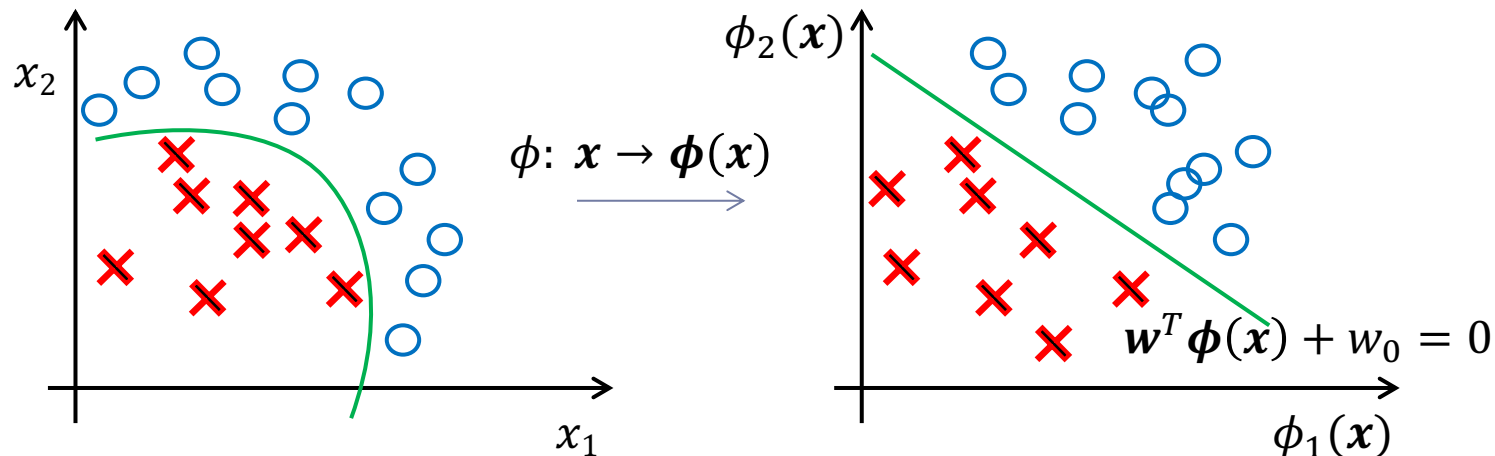
$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{m-1}x^{m-1} + w_mx^m$$

# Polynomial regression: example



# Classification: Not linearly separable data

- Non-linear decision surface: Transform to a new feature space



- Quadratic surfaces
  - Two dimensional feature space:

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]^T$$

- d-dimensional feature space

$$\phi(\mathbf{x}) = [1, x_1, \dots, x_d, x_1^2, \dots, x_d^2, x_1x_2, \dots, x_1x_d, x_2x_3, \dots, x_{d-1}x_d]^T$$



# Model complexity and overfitting

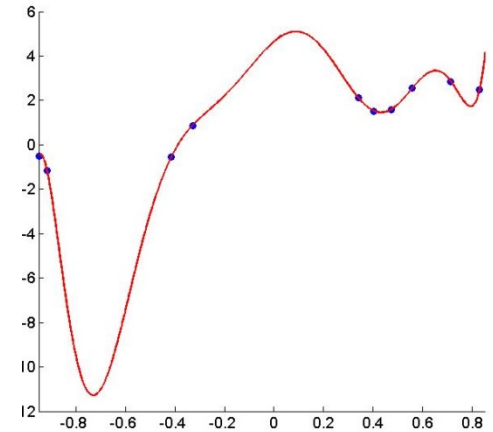
- With limited training data, models may achieve zero training error but a large test error.

Training  
(empirical) loss

$$\frac{1}{n} \sum_{i=1}^n \left( y^{(i)} - f(\mathbf{x}^{(i)}; \boldsymbol{\theta}) \right)^2 \approx 0$$

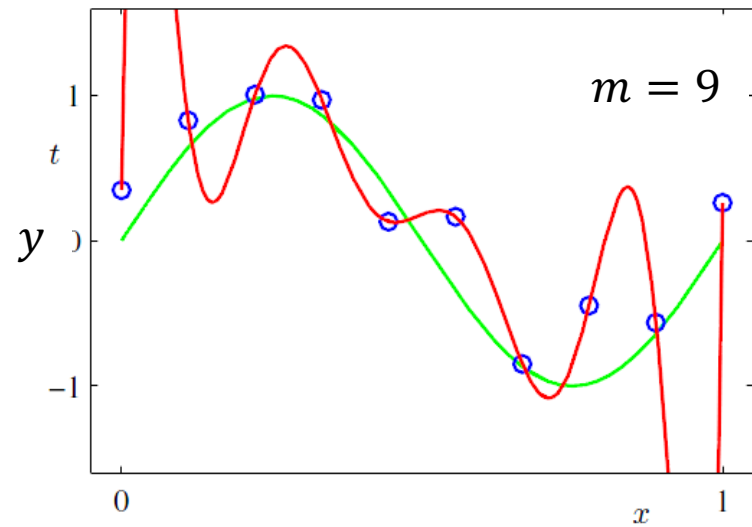
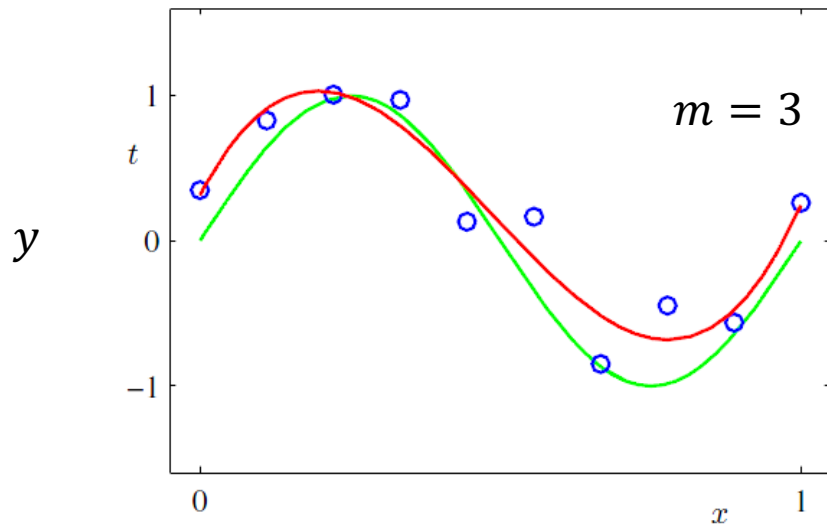
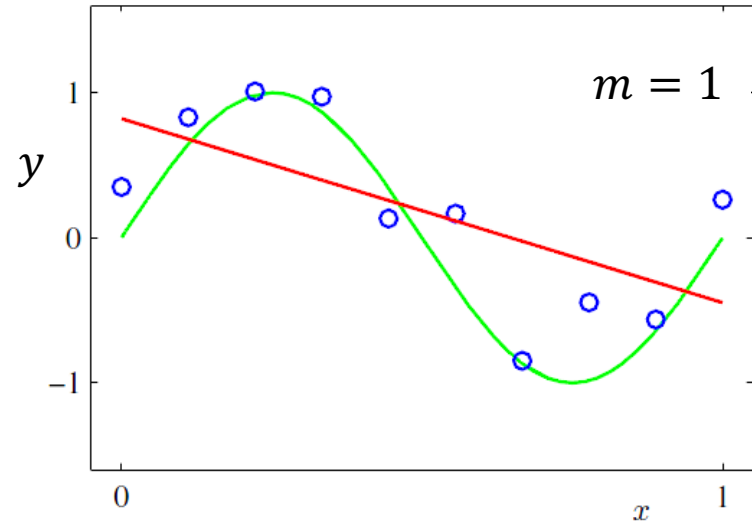
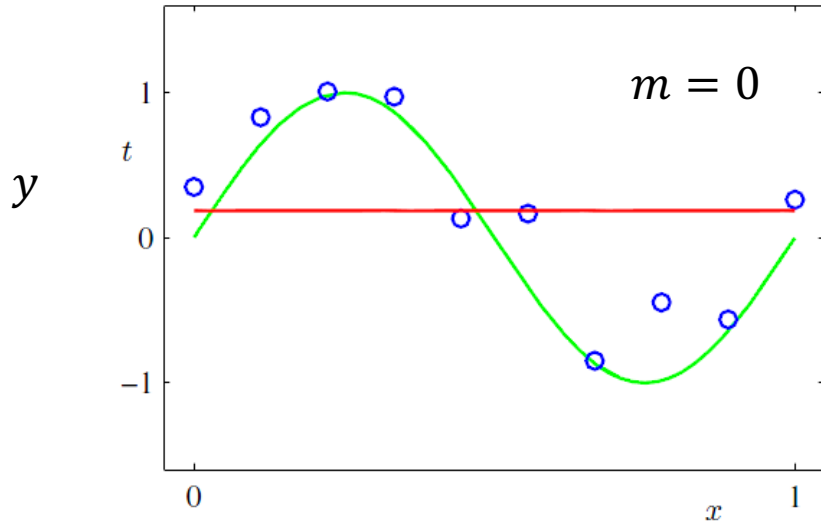
Expected  
(true) loss

$$E_{\mathbf{x}, y} \left\{ (y - f(\mathbf{x}; \boldsymbol{\theta}))^2 \right\} \gg 0$$



- Over-fitting: when the training loss no longer bears any relation to the test (generalization) loss.
  - Fails to generalize to unseen examples.

# Polynomial regression

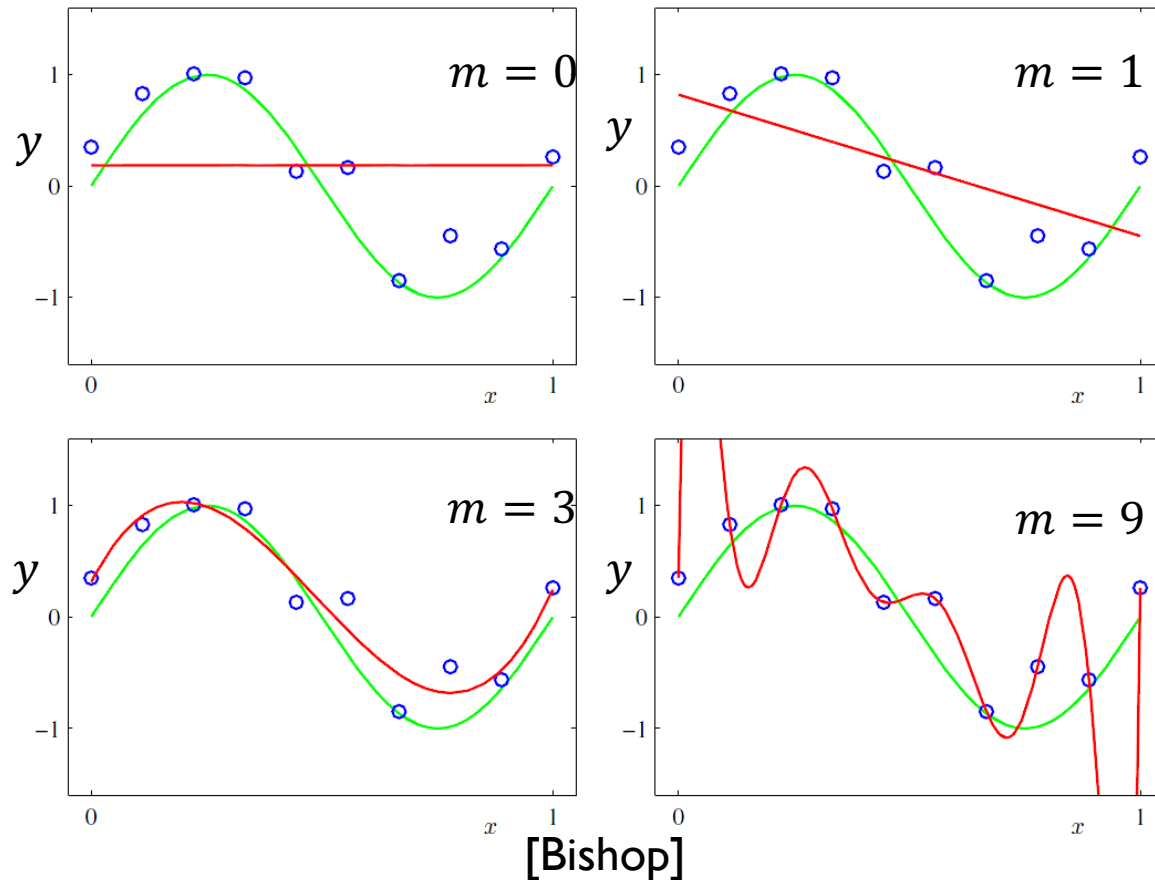


# Over-fitting causes

- **Model complexity**
  - E.g., Model with a large number of parameters (degrees of freedom)
- **Low number of training data**
  - Small data size compared to the complexity of the model

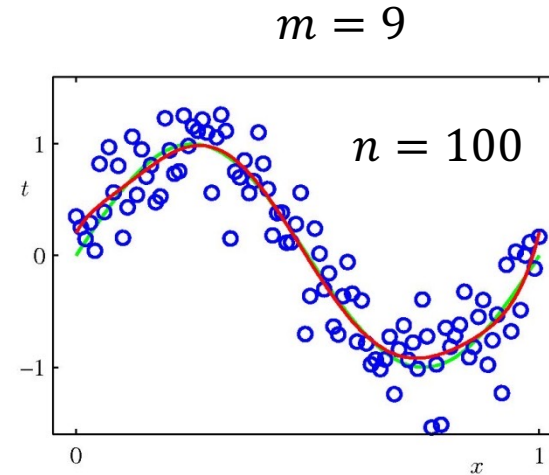
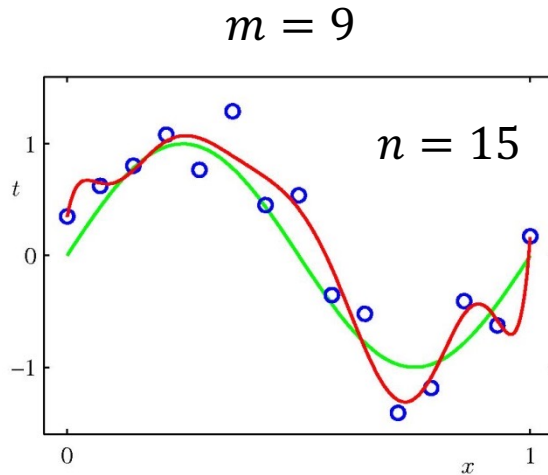
# Model complexity

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



# Number of training data & overfitting

- ▶ Over-fitting problem becomes less severe as the size of training data increases.



[Bishop]

# Avoiding over-fitting

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

# Avoiding over-fitting

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

# Evaluation and model selection

- **Evaluation:**
  - We need to measure how well the learned function can predict 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



# Model selection

- **Learning algorithm** defines the data-driven search over the hypothesis space
  - Optimization of parameters
- **Hyper-parameters** are the tunable aspects of the model, that the learning algorithm does *not* select

# Model selection

- **Model selection** is the process by which we choose the “best” model among a set of candidates
  - assume access to a function capable of measuring the quality of a model
  - typically done “outside” the main training algorithm
- Model selection / hyper-parameter optimization is just another form of learning

# Simple hold-out: model selection

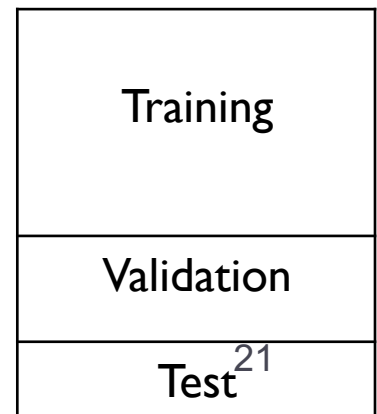
- Steps:
  - Divide training data into training and validation set  $v\_set$
  - Use only the training set to train a set of models
  - Evaluate each learned model on the validation set
    - $J_v(\mathbf{w}) = \frac{1}{|v\_set|} \sum_{i \in v\_set} \left( y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$
  - Choose the best model based on the validation set error

# 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 obtains a relatively noisy estimate of performance.

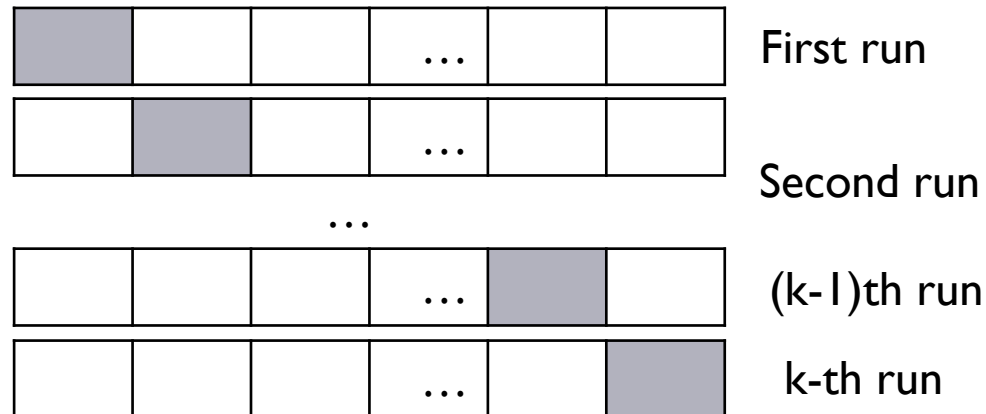
# Simple hold out: training, validation, and test sets

- Simple hold-out chooses the model that minimizes error on validation set.
- $J_v(\hat{\mathbf{W}})$  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



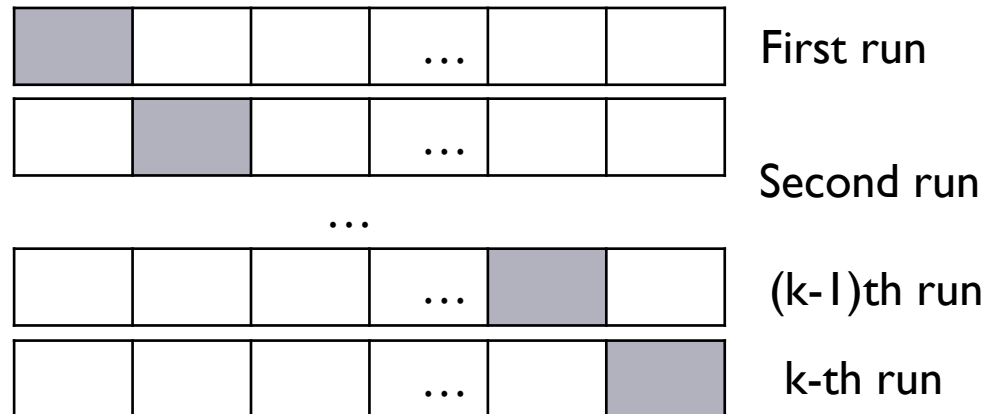
# Cross-Validation (CV): evaluation

- $k$ -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



# 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 as an estimation of the true performance of the model.



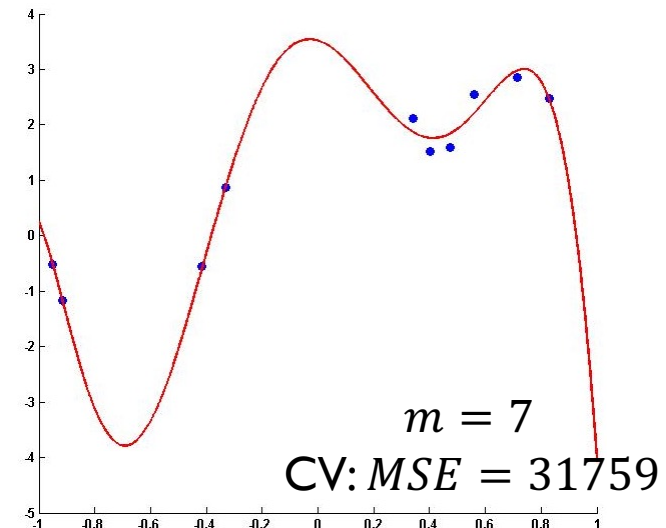
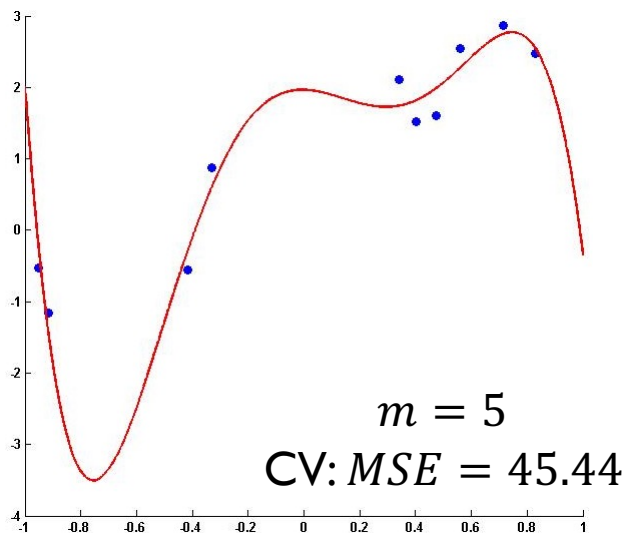
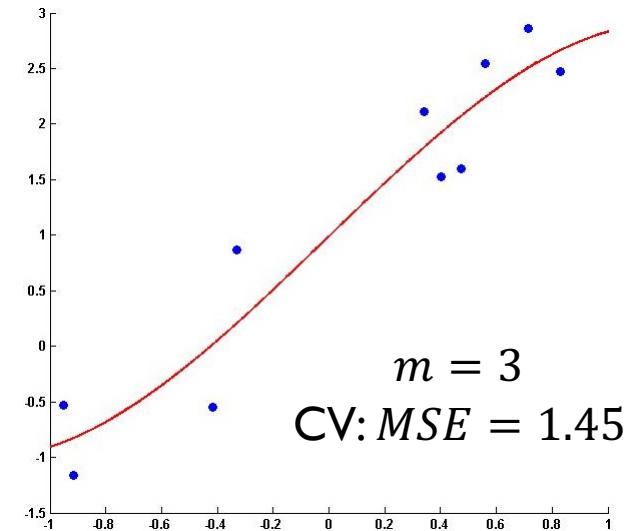
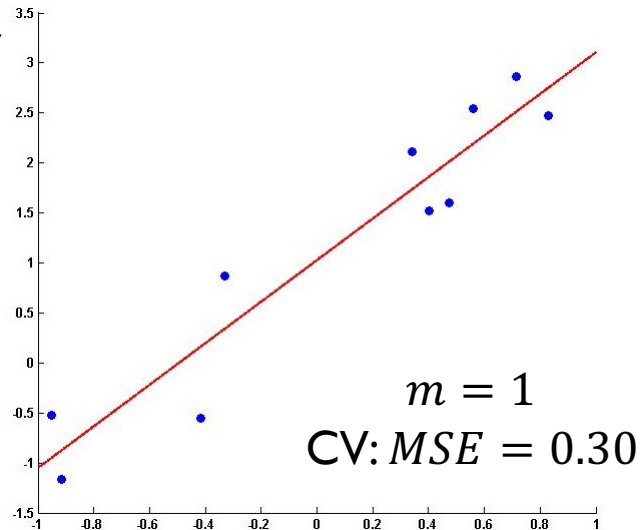
# Cross-Validation (CV): model selection

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



# Cross-validation: polynomial regression example

- 5-fold CV
- 100 runs
- average



# Avoiding over-fitting

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

# Regularization

- Adding a penalty term in the cost function to discourage the coefficients from reaching large values.

# Regularization in regression problem

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

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

# Regularization in regression problem

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

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

- Weight update by gradient descent:

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

# Regularization in classification problem

---

- Multi-class logistic regression (i.e., cross entropy loss) with regularization:

$$J(\mathbf{W}) = - \sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log \left( g_k(\mathbf{x}^{(i)}; \mathbf{W}) \right) + \lambda \sum_{k=1}^K \mathbf{w}_k^T \mathbf{w}_k$$

- Weight Update:

$$\mathbf{w}_k^{t+1} = \mathbf{w}_k^t - \eta \nabla_{\mathbf{w}_k} J(\mathbf{W}^t)$$
$$\nabla_{\mathbf{w}_k} J(\mathbf{W}) = -2 \sum_{i=1}^n (y^{(i)} - g_k(\mathbf{x}^{(i)}; \mathbf{W})) \mathbf{x}^{(i)} + 2\lambda \mathbf{w}_k$$

# Regression: 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  $m$ .

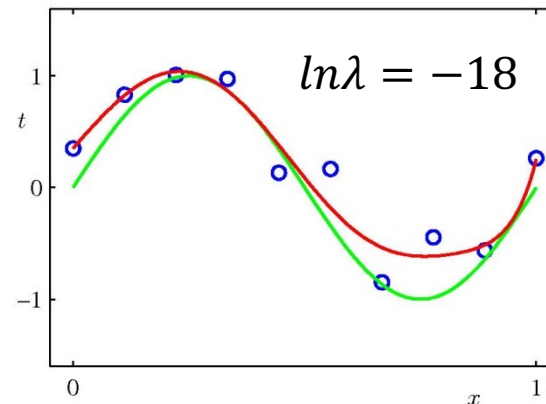
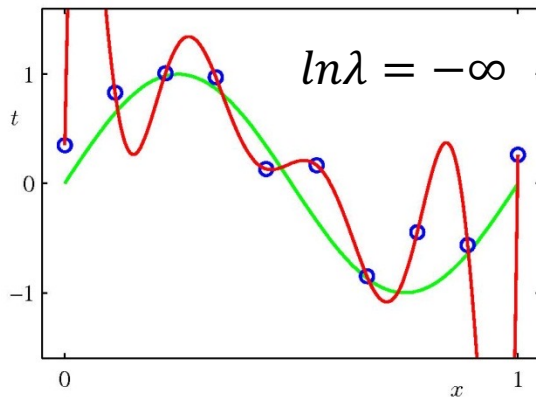
	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

[Bishop]

# Regression: regularization parameter

	$m = 9$		
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$\hat{w}_0$	0.35	0.35	0.13
$\hat{w}_1$	232.37	4.74	-0.05
$\hat{w}_2$	-5321.83	-0.77	-0.06
$\hat{w}_3$	48568.31	-31.97	-0.05
$\hat{w}_4$	-231639.30	-3.89	-0.03
$\hat{w}_5$	640042.26	55.28	-0.02
$\hat{w}_6$	-1061800.52	41.32	-0.01
$\hat{w}_7$	1042400.18	-45.95	-0.00
$\hat{w}_8$	-557682.99	-91.53	0.00
$\hat{w}_9$	125201.43	72.68	0.01

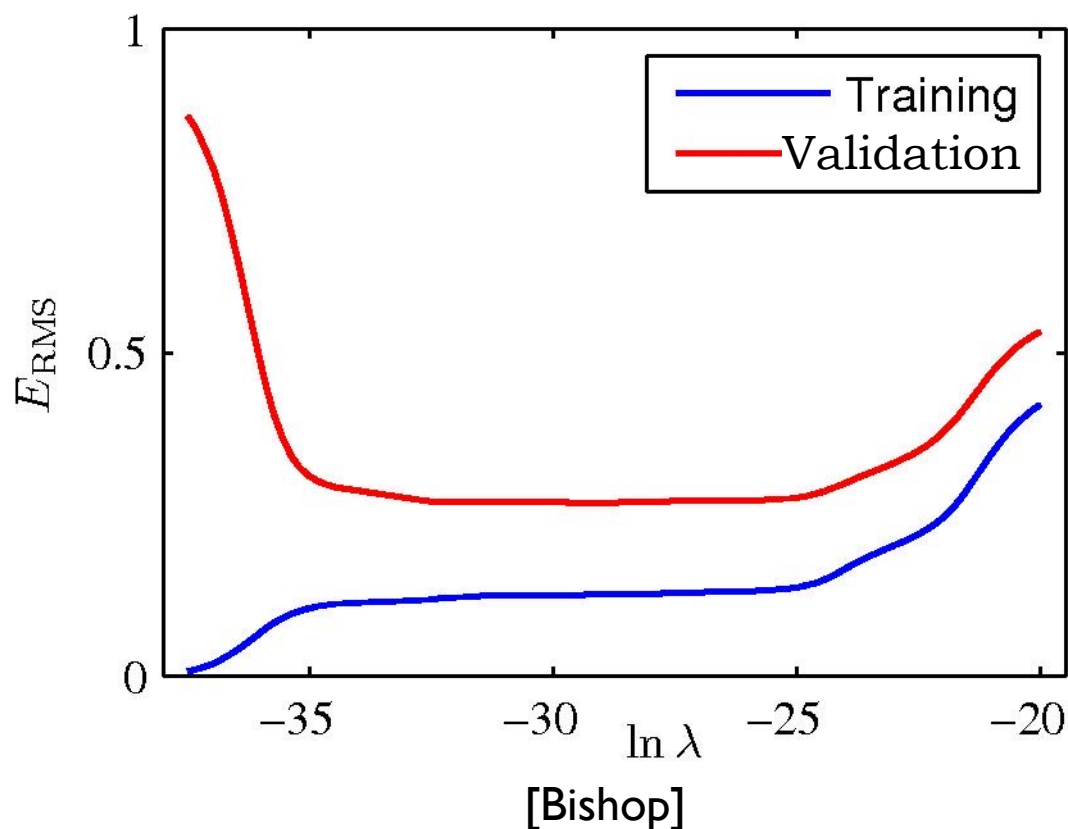
[Bishop]





# Regularization parameter

- Generalization
  - $\lambda$  now controls the effective complexity of the model and hence determines the degree of over-fitting



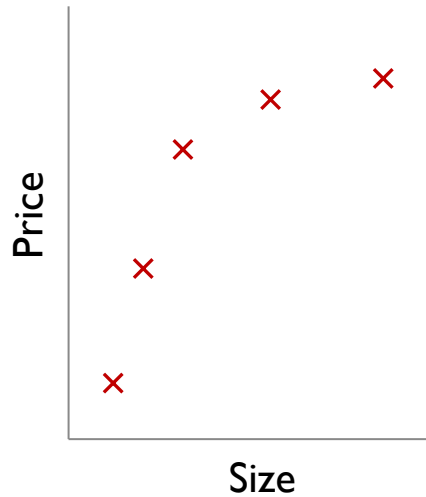
# Choosing the regularization parameter

- A set of models with different values of  $\lambda$ .
- Find  $\hat{\mathbf{W}}$  for each model based on training data
- Find  $J_v(\hat{\mathbf{W}})$  (or  $J_{cv}(\hat{\mathbf{W}})$ ) for each model
  - $J_v(\mathbf{w}) = \frac{1}{n_v} \sum_{i \in v\_set} \left( y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$
- Select the model with the best  $J_v(\hat{\mathbf{W}})$  (or  $J_{cv}(\hat{\mathbf{W}})$ )

# 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

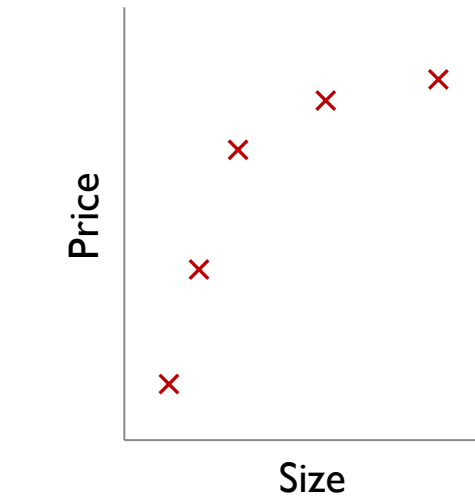


$$w_0 + w_1x$$

Less complex  $\mathcal{H}$



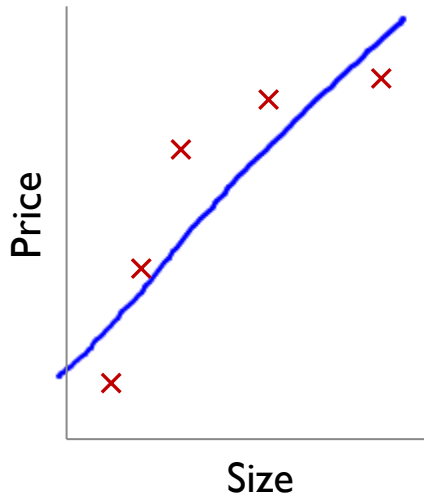
$$w_0 + w_1x + w_2x^2$$



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

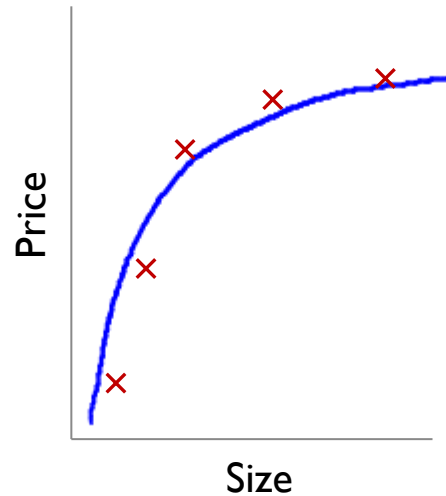
More complex  $\mathcal{H}$

# Complexity of hypothesis space: example

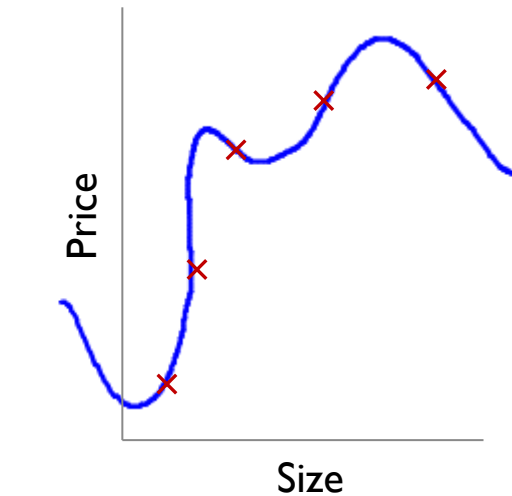


$$w_0 + w_1x$$

Underfitting



$$w_0 + w_1x + w_2x^2$$



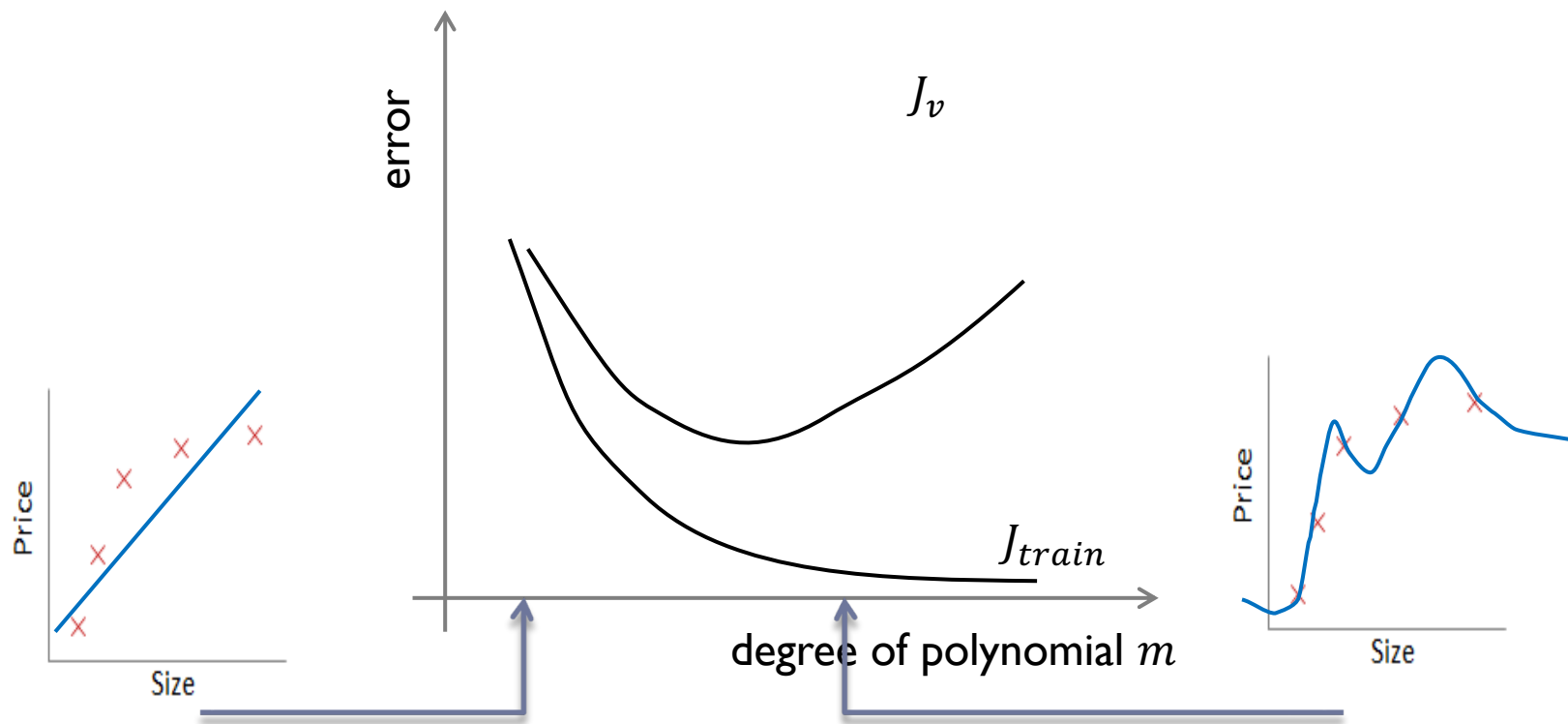
$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

Overfitting

# Complexity of hypothesis space: example

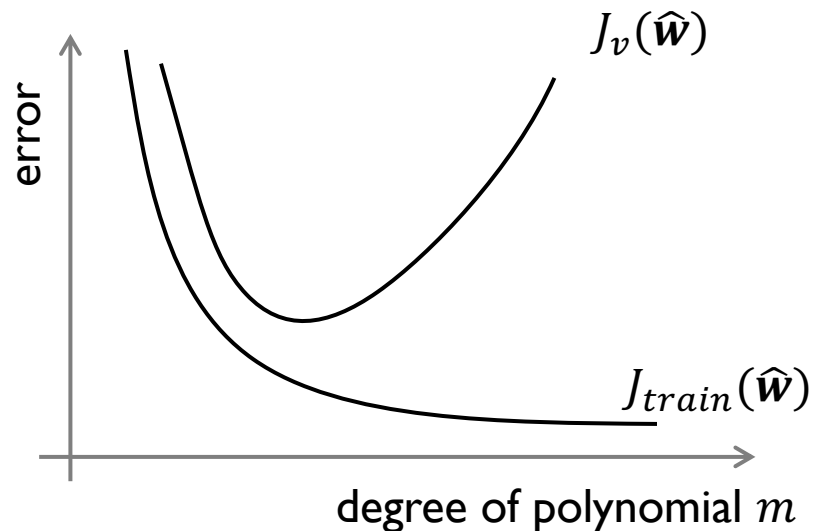
$$J_v(\mathbf{w}) = \frac{1}{n_v} \sum_{i \in \text{eval\_set}} \left( y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$

$$J_{\text{train}}(\mathbf{w}) = \frac{1}{n_{\text{train}}} \sum_{i \in \text{train\_set}} \left( y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$



# Complexity of hypothesis space

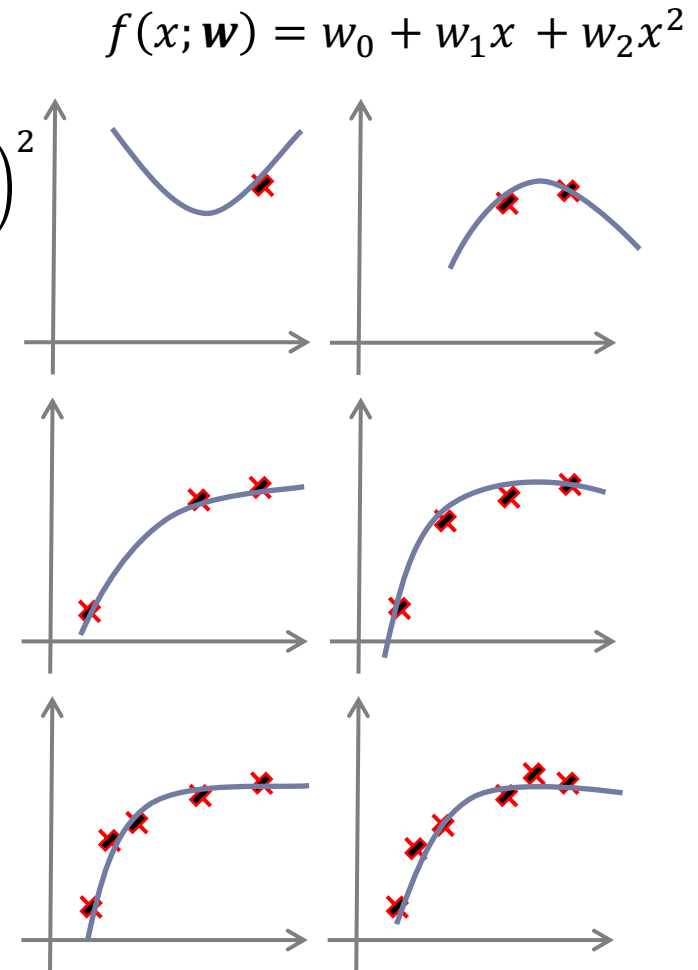
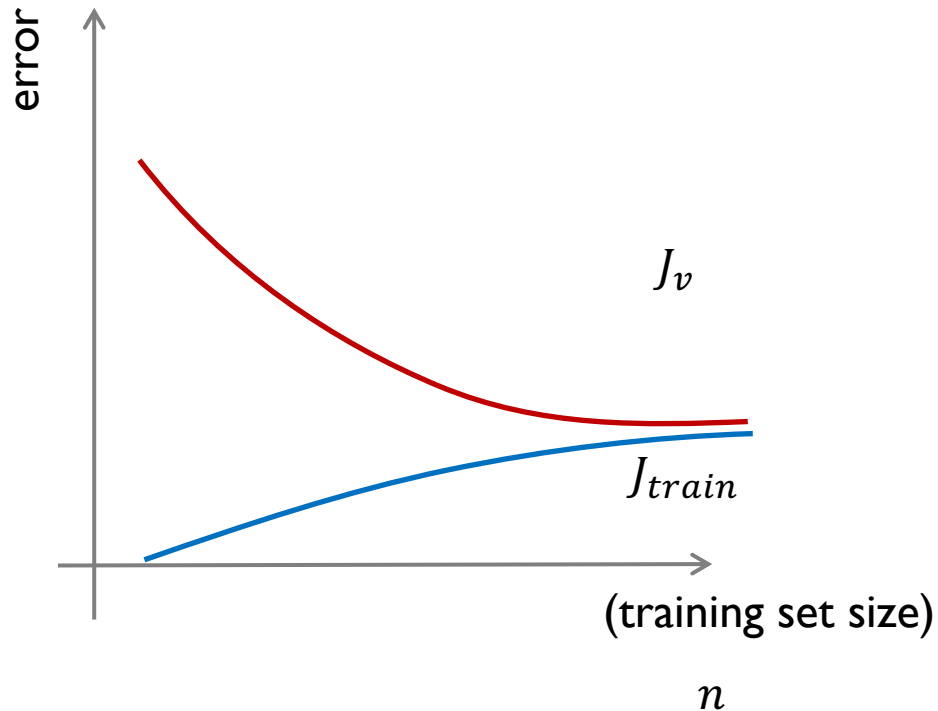
- Less complex  $\mathcal{H}$ :
  - $J_{train}(\hat{\mathbf{w}}) \approx J_v(\hat{\mathbf{w}})$  and  $J_{train}(\hat{\mathbf{w}})$  is very high
- More complex  $\mathcal{H}$ :
  - $J_{train}(\hat{\mathbf{w}}) \ll J_v(\hat{\mathbf{w}})$  and  $J_{train}(\hat{\mathbf{w}})$  is low



# Size of training set

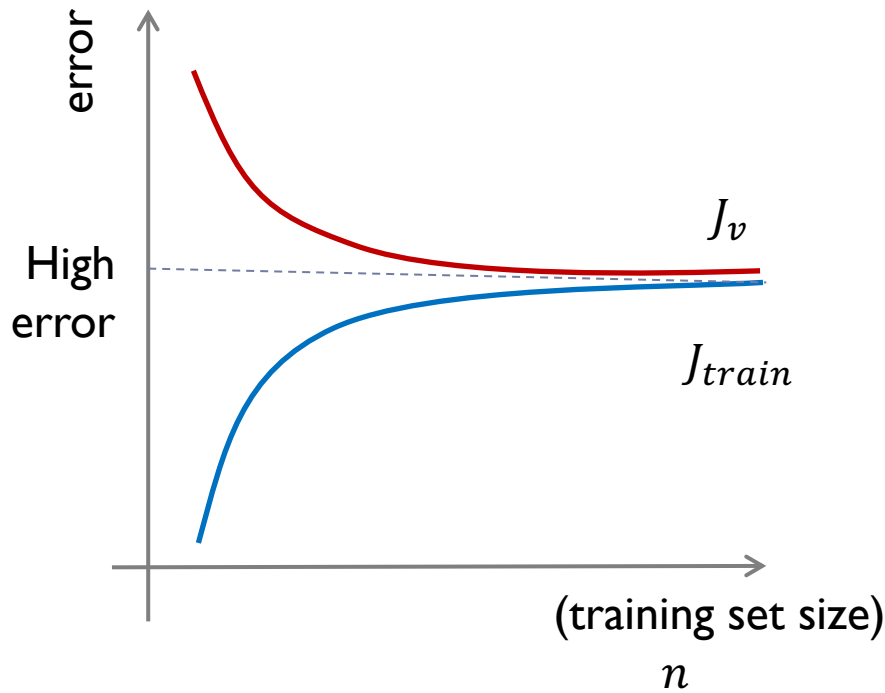
$$J_v(\mathbf{w}) = \frac{1}{n_{val}} \sum_{i \in val\_set} \left( y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$

$$J_{train}(\mathbf{w}) = \frac{1}{n_{train}} \sum_{i \in train\_set} \left( y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$

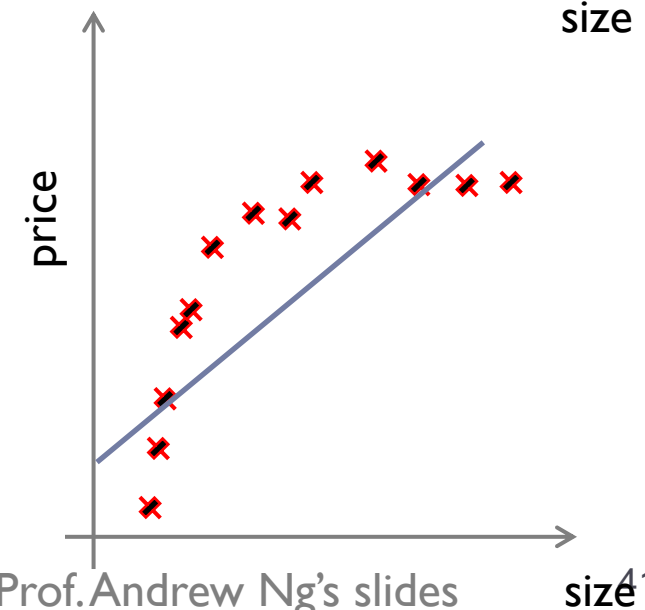
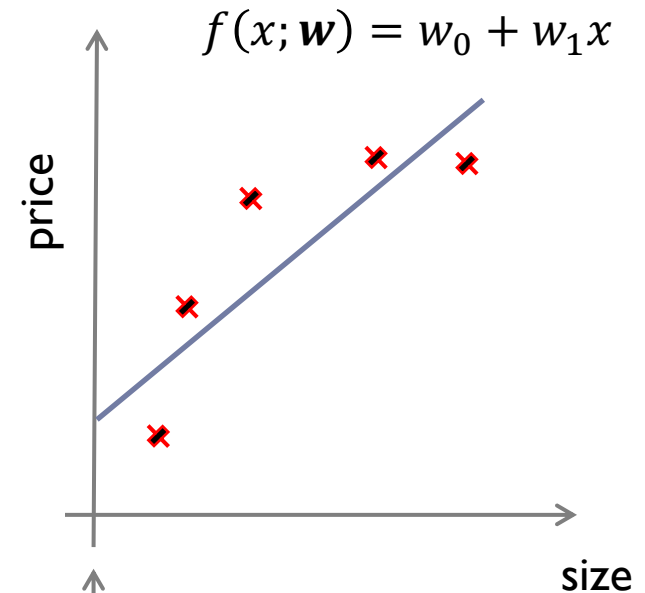




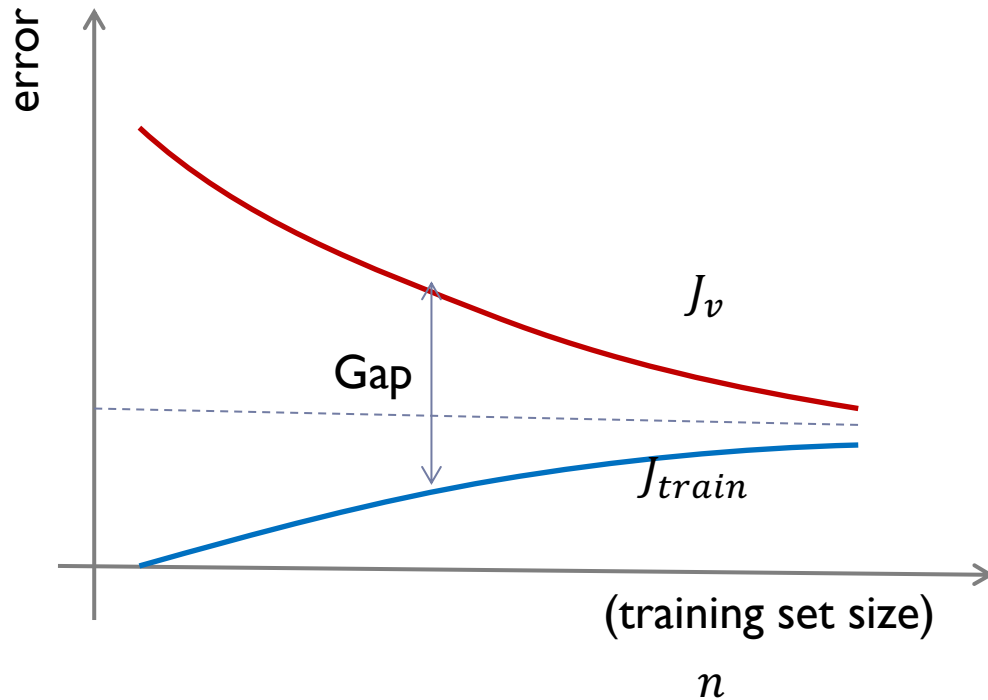
# Less complex $\mathcal{H}$



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

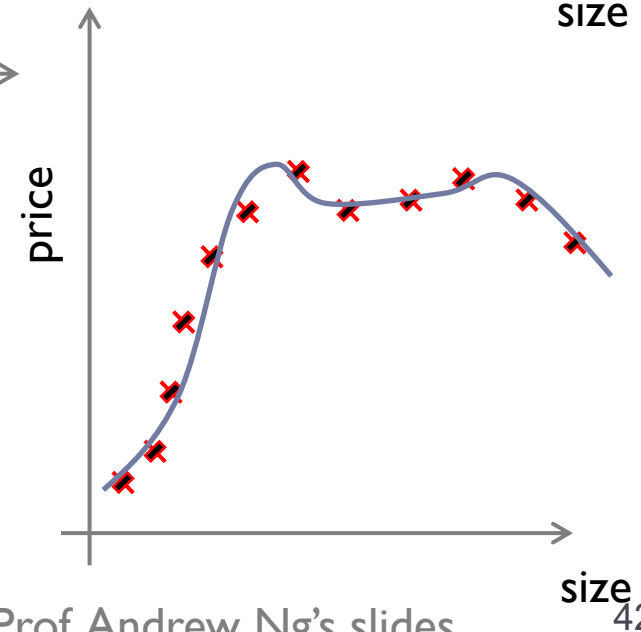
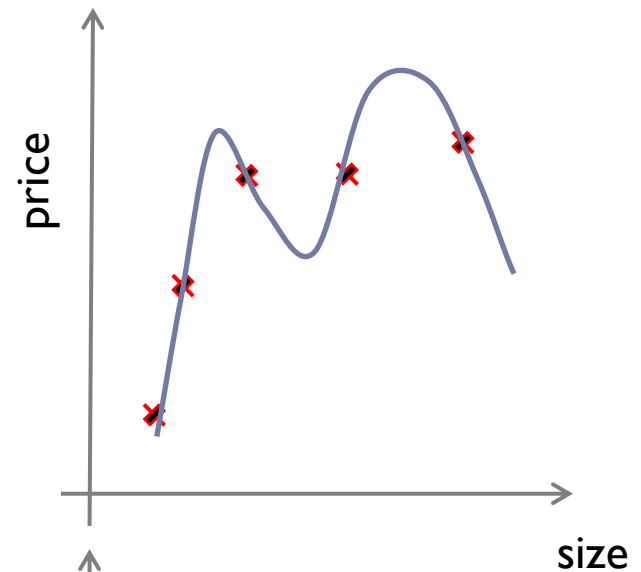


# More complex $\mathcal{H}$



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

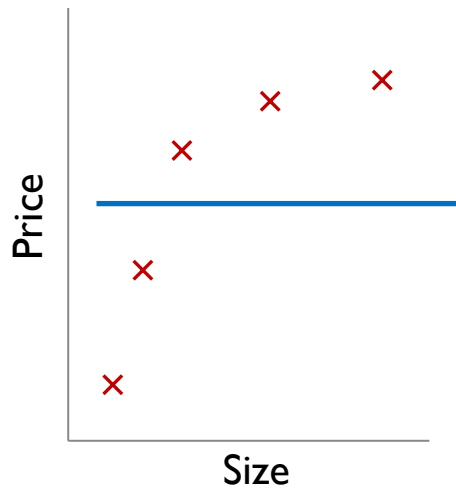
$$f(x; \mathbf{w}) = w_0 + w_1x + \dots + w_{10}x^{10}$$



# Regularization: example

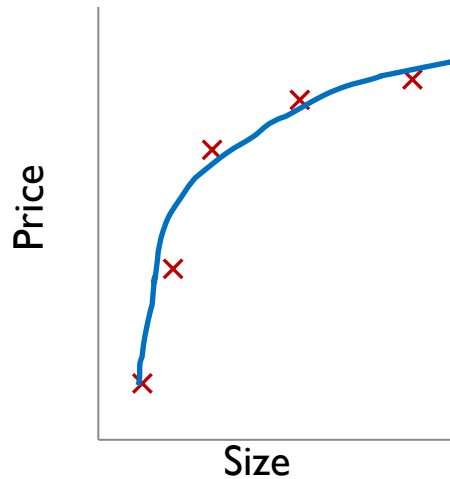
$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

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

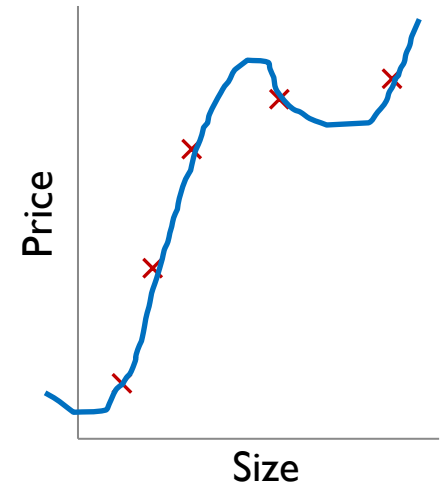


Large  $\lambda$   
(Prefer to more simple models)

$$w_1 = w_2 \approx 0$$



Intermediate  $\lambda$



Small  $\lambda$   
(Prefer to more complex models)

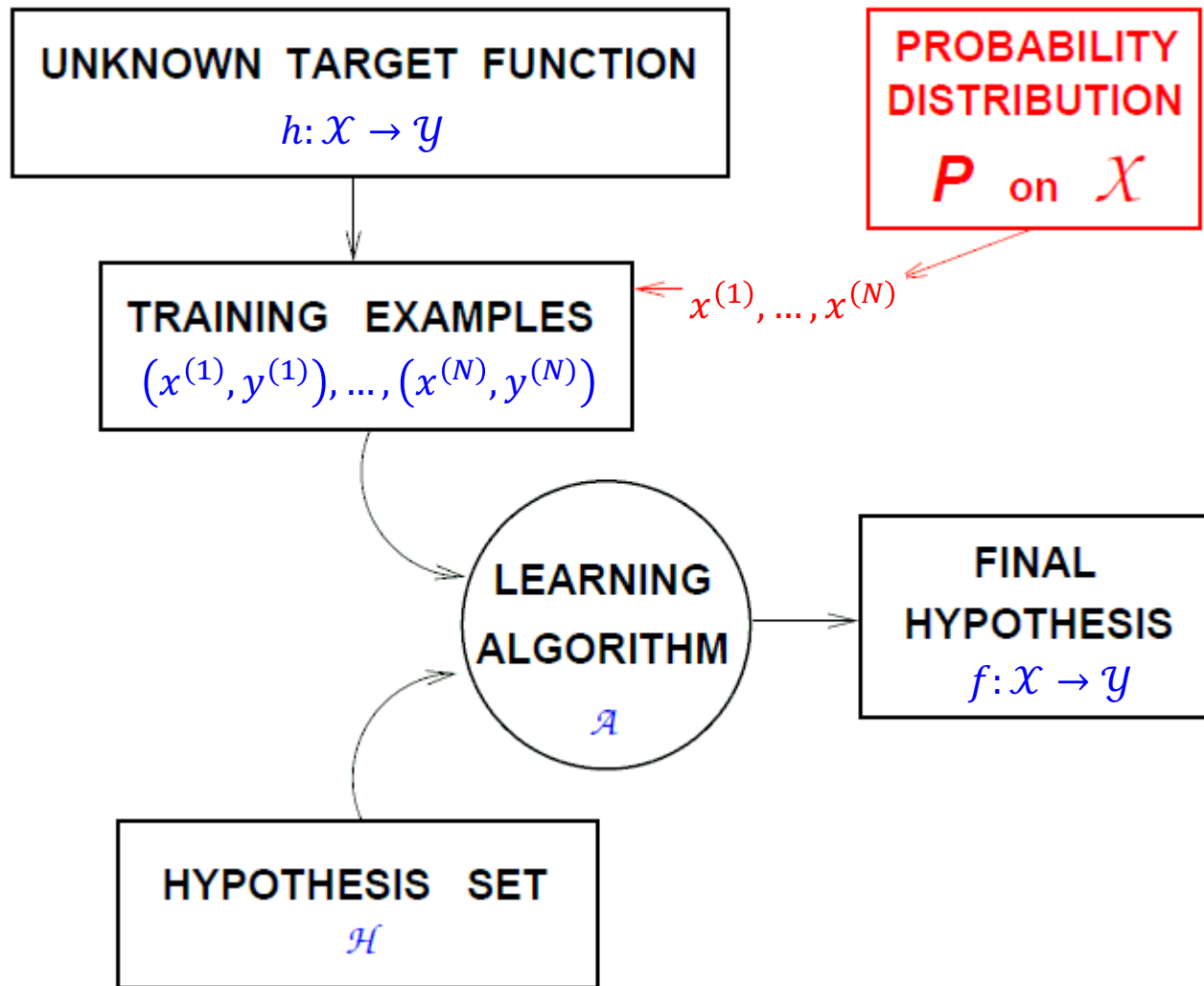
$$\lambda = 0$$

# Model complexity: Bias-variance trade-off

---

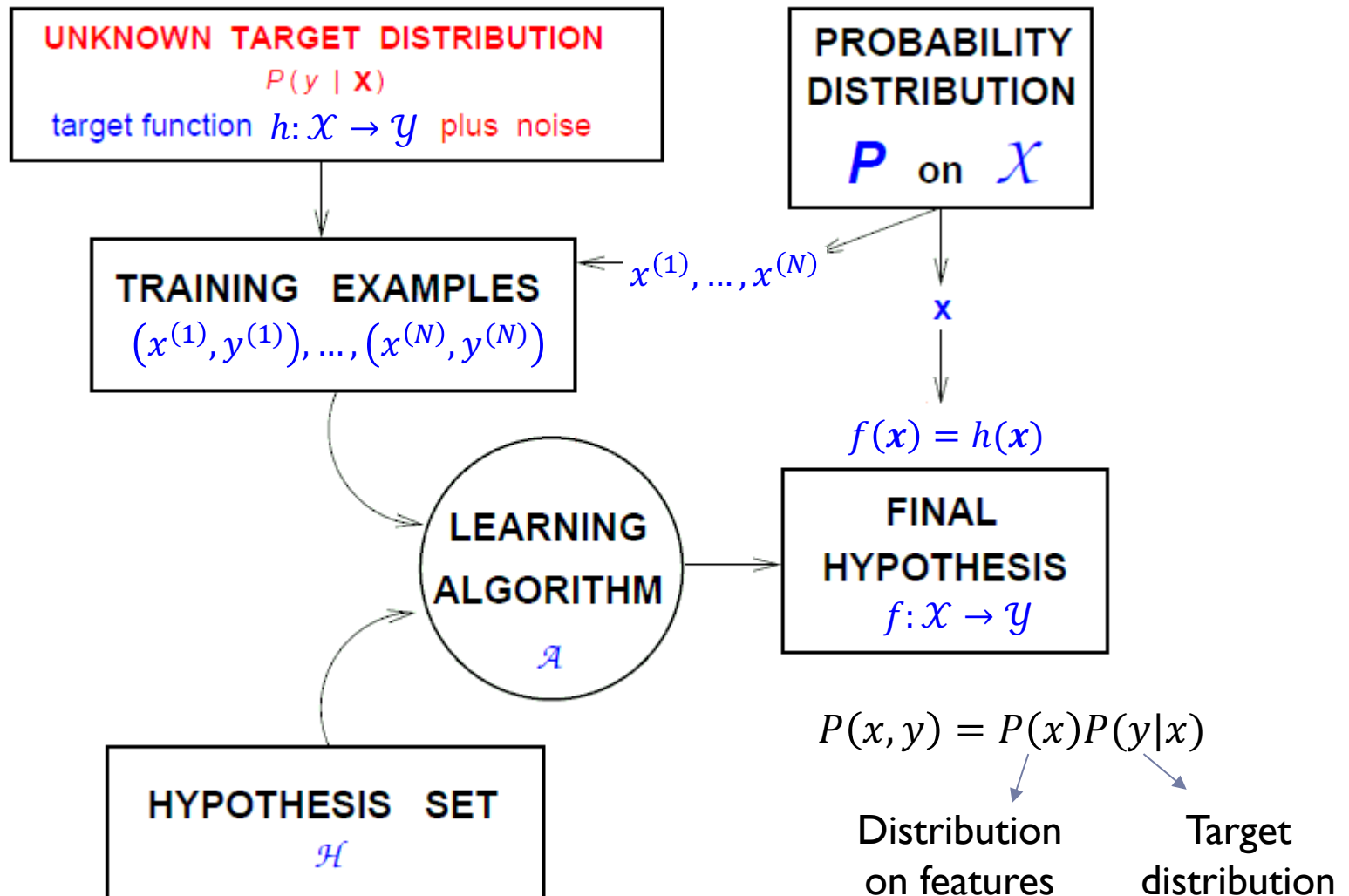
- Least squares can lead to severe over-fitting if complex models are trained using data sets of limited size.
- A frequentist viewpoint of the model complexity issue, known as the *bias-variance trade-off*.

# The learning diagram: deterministic target



# The learning diagram including noisy target

- Type



# Expectation of true error $(\mathbf{x}, y) \sim P$

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

---

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

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

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

---

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



## Using the average hypothesis

---

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

## Using the average hypothesis

---

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

# Bias and variance

---

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

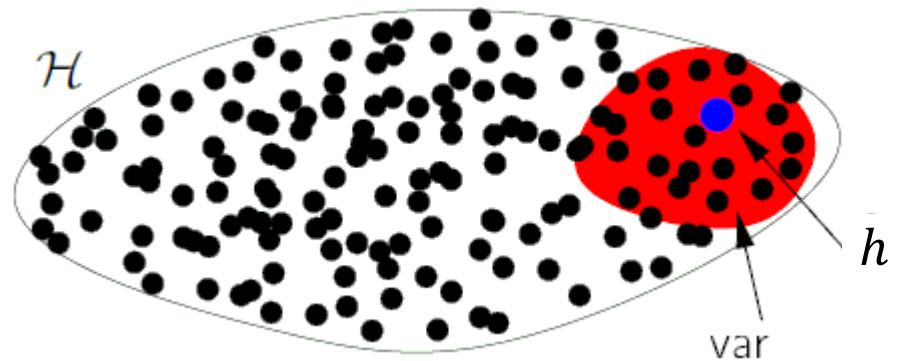
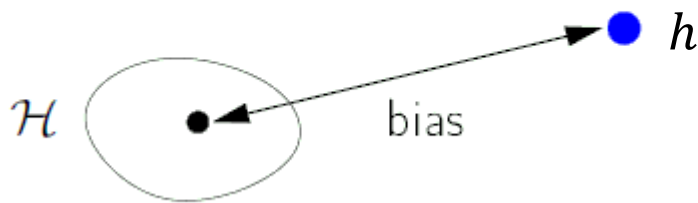
$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathcal{D}} \left[ (f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] \right] &= \mathbb{E}_{\mathbf{x}} [\text{var}(\mathbf{x}) + \text{bias}(\mathbf{x})] \\ &= \text{var} + \text{bias} \end{aligned}$$

# Bias-variance trade-off

---

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

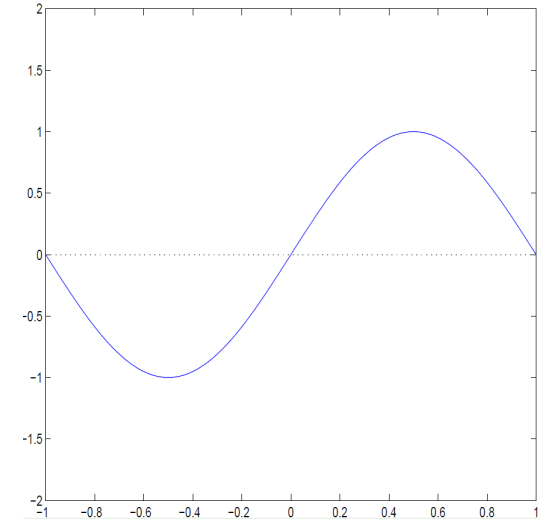
$$\text{bias} = \mathbb{E}_{\mathbf{x}} \left[ \bar{f}(\mathbf{x}) - h(\mathbf{x}) \right]$$



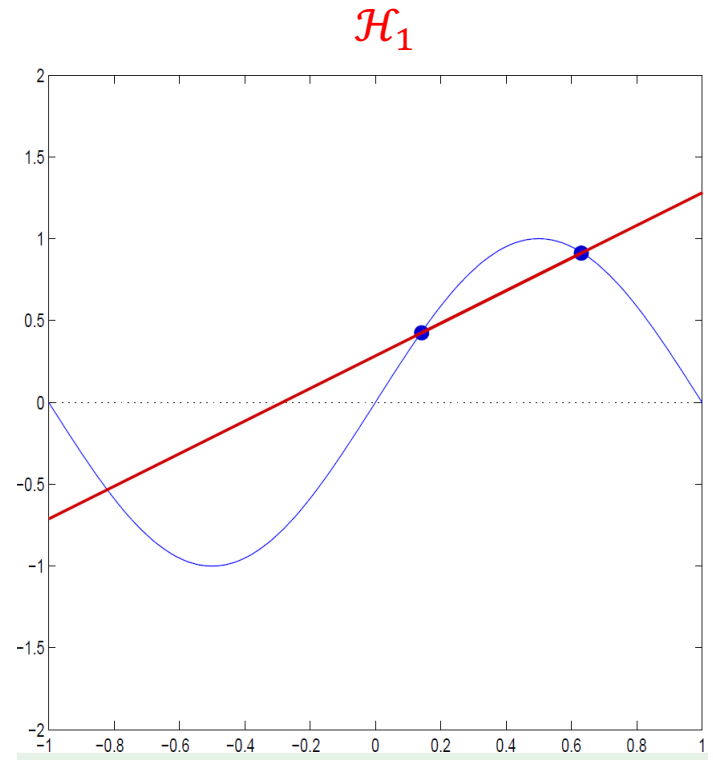
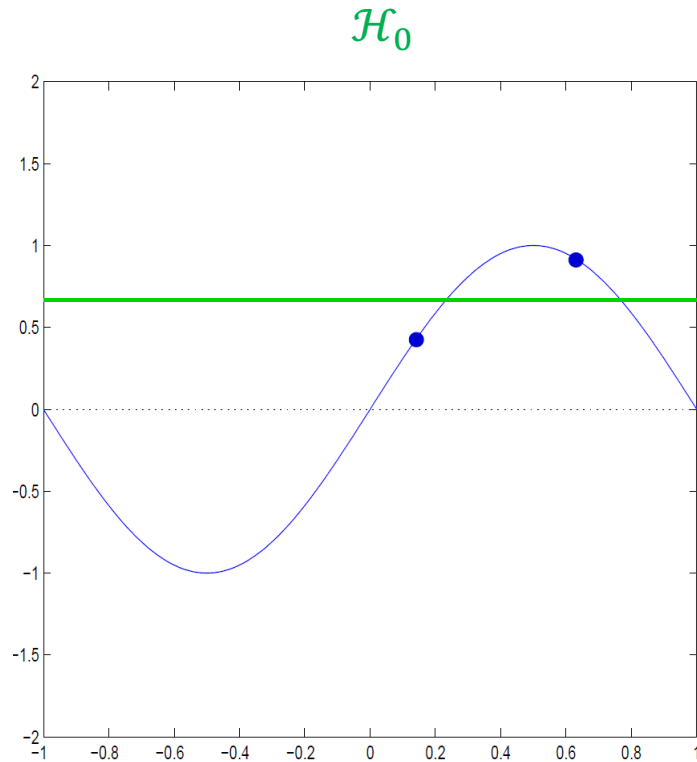
More complex  $\mathcal{H} \Rightarrow$  lower bias but higher variance

# Example I: sin target

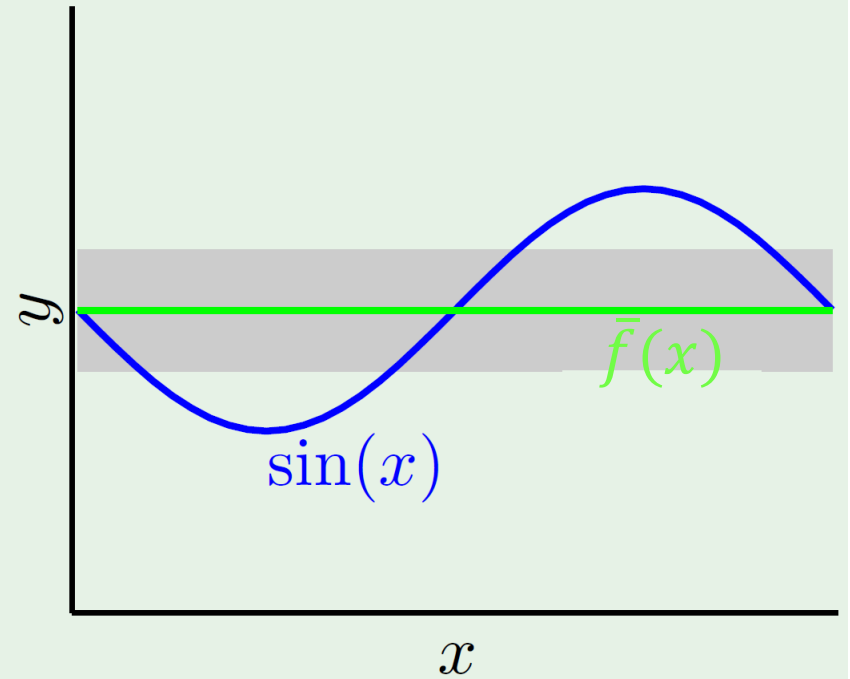
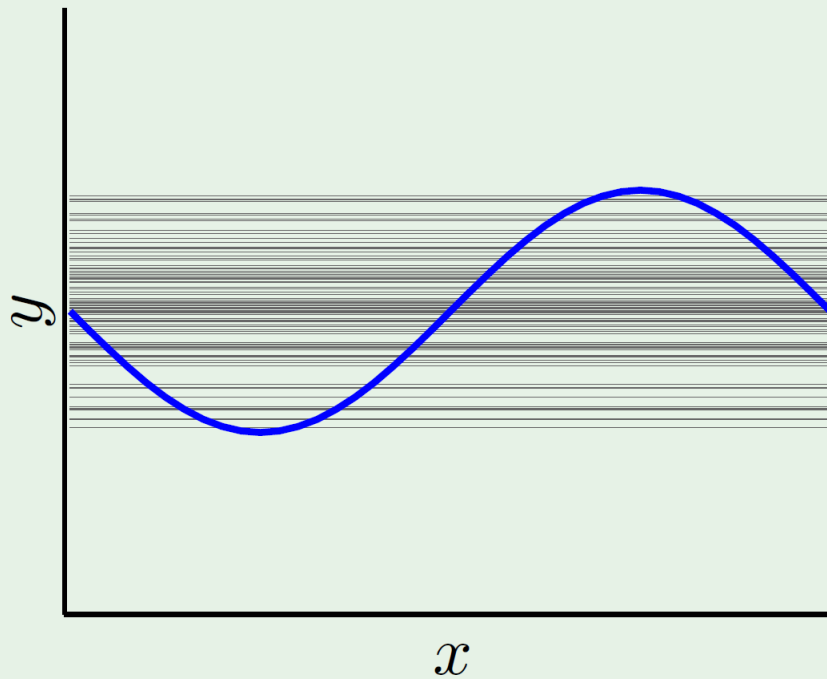
- Only two training example  $N = 2$
- Two models used for learning:
  - $\mathcal{H}_0: f(x) = b$
  - $\mathcal{H}_1: f(x) = ax + b$
- Which is better  $\mathcal{H}_0$  or  $\mathcal{H}_1$ ?



# Example I: learning from a training set

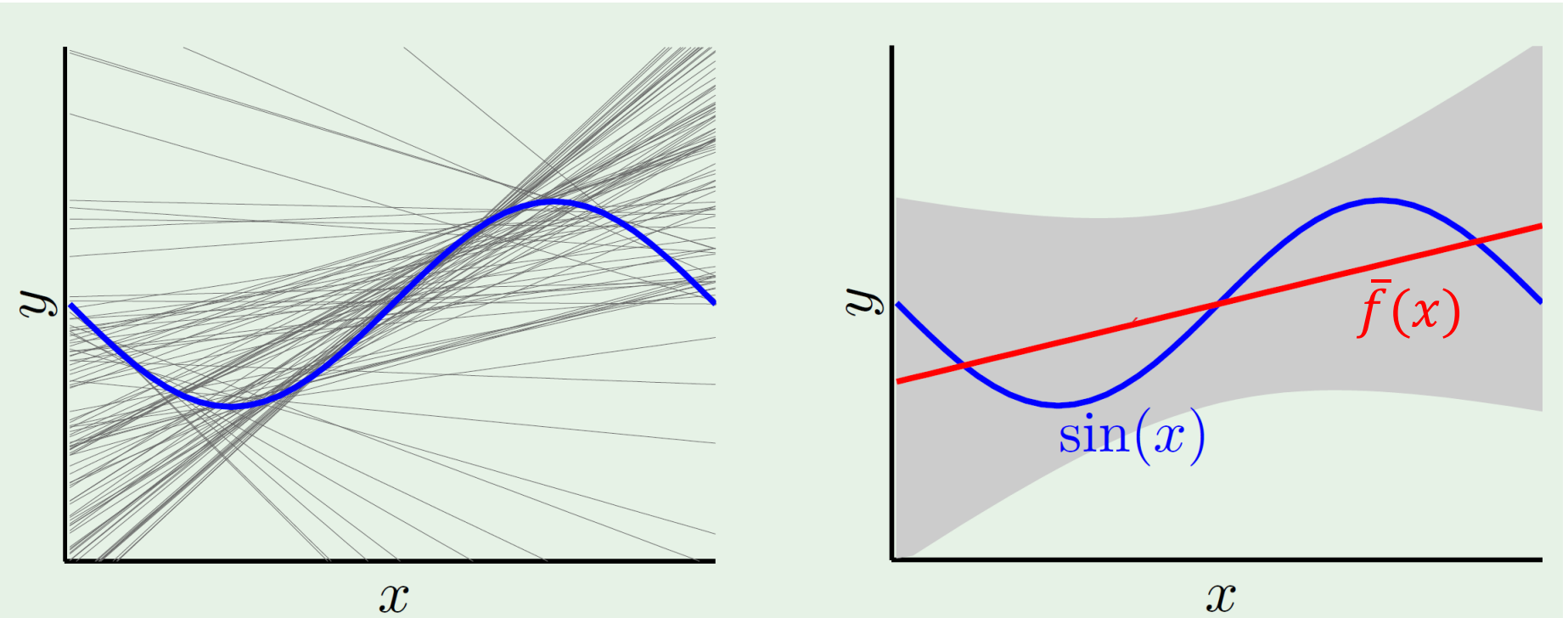


# Example I: variance $\mathcal{H}_0$

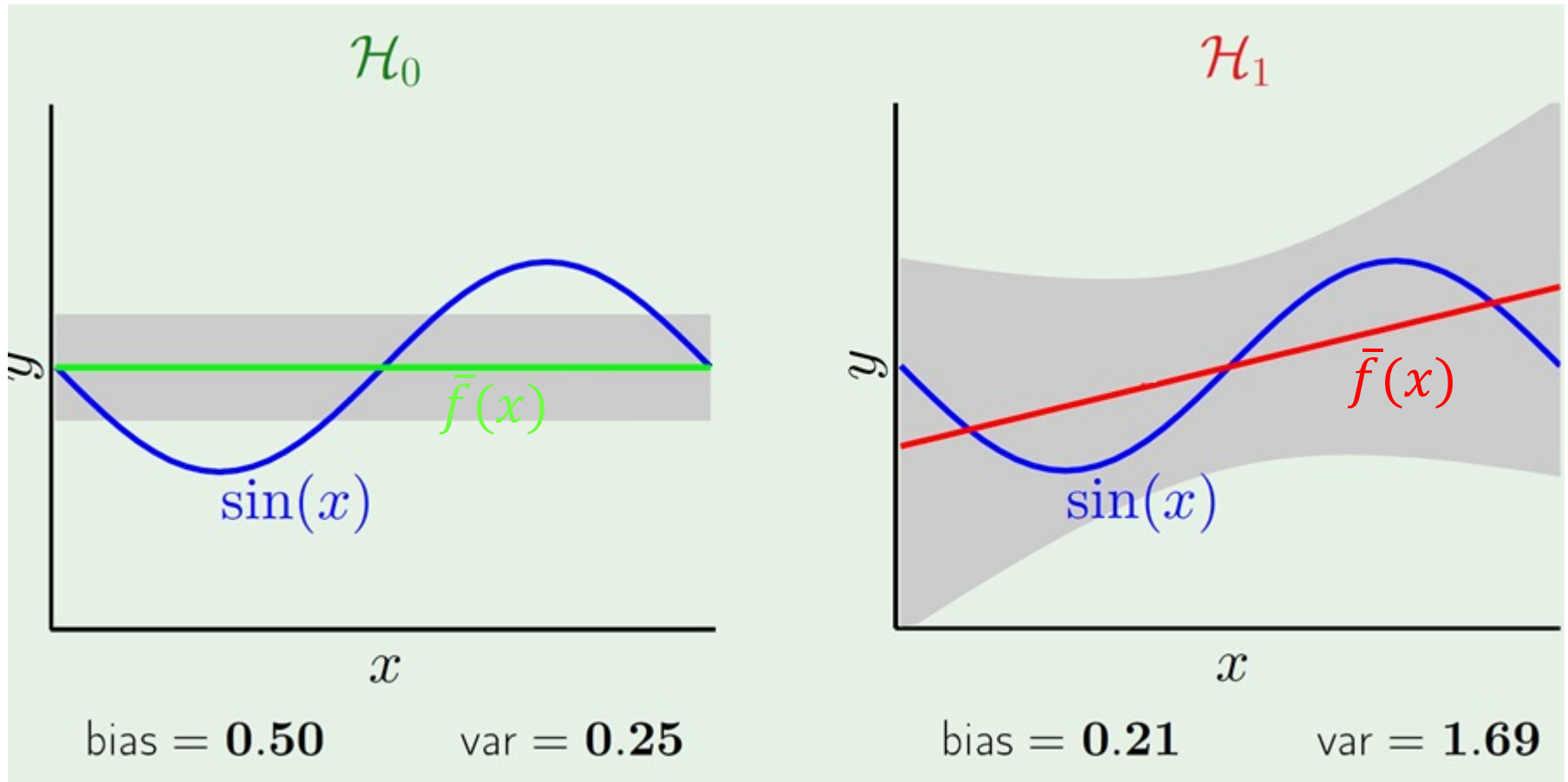




# Example I: variance $\mathcal{H}_1$



# Example I: which is better?



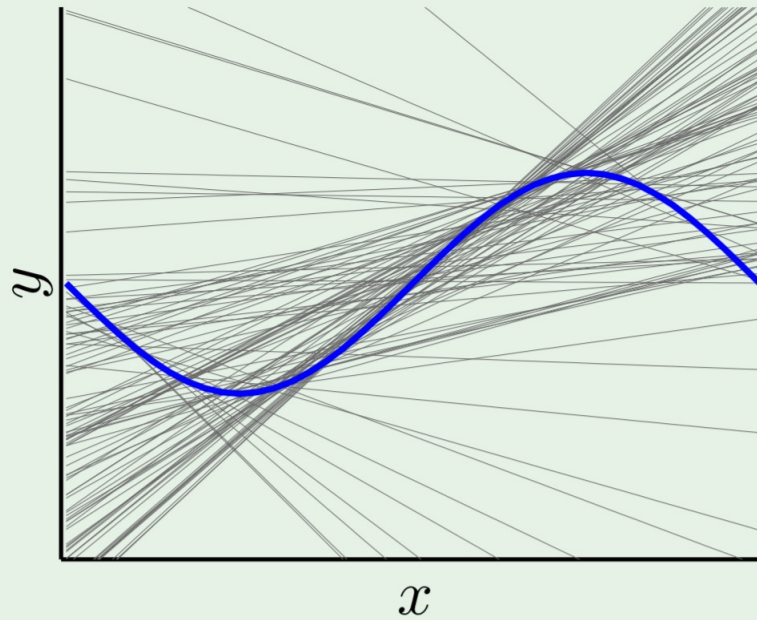
# Lesson

Match the **model complexity**

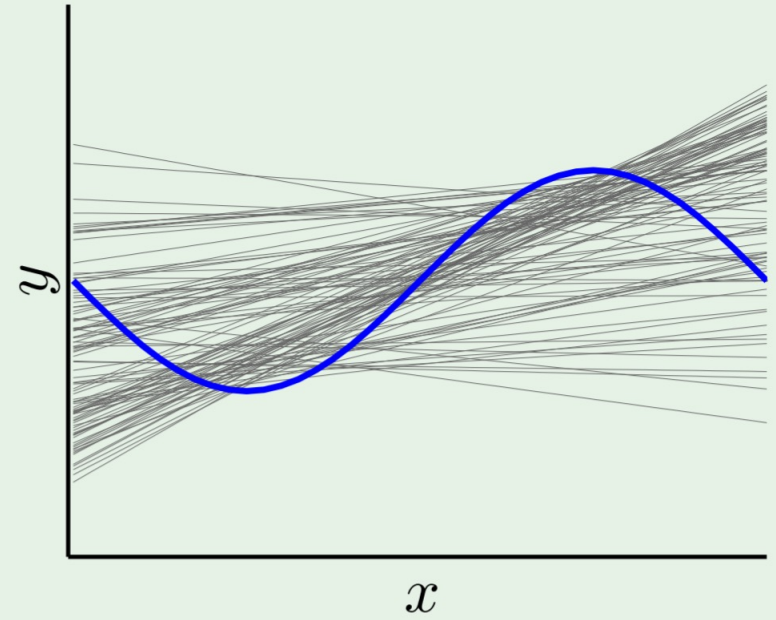
to the **data sources**

not to the complexity of the **target function**.

# Example I: regularization

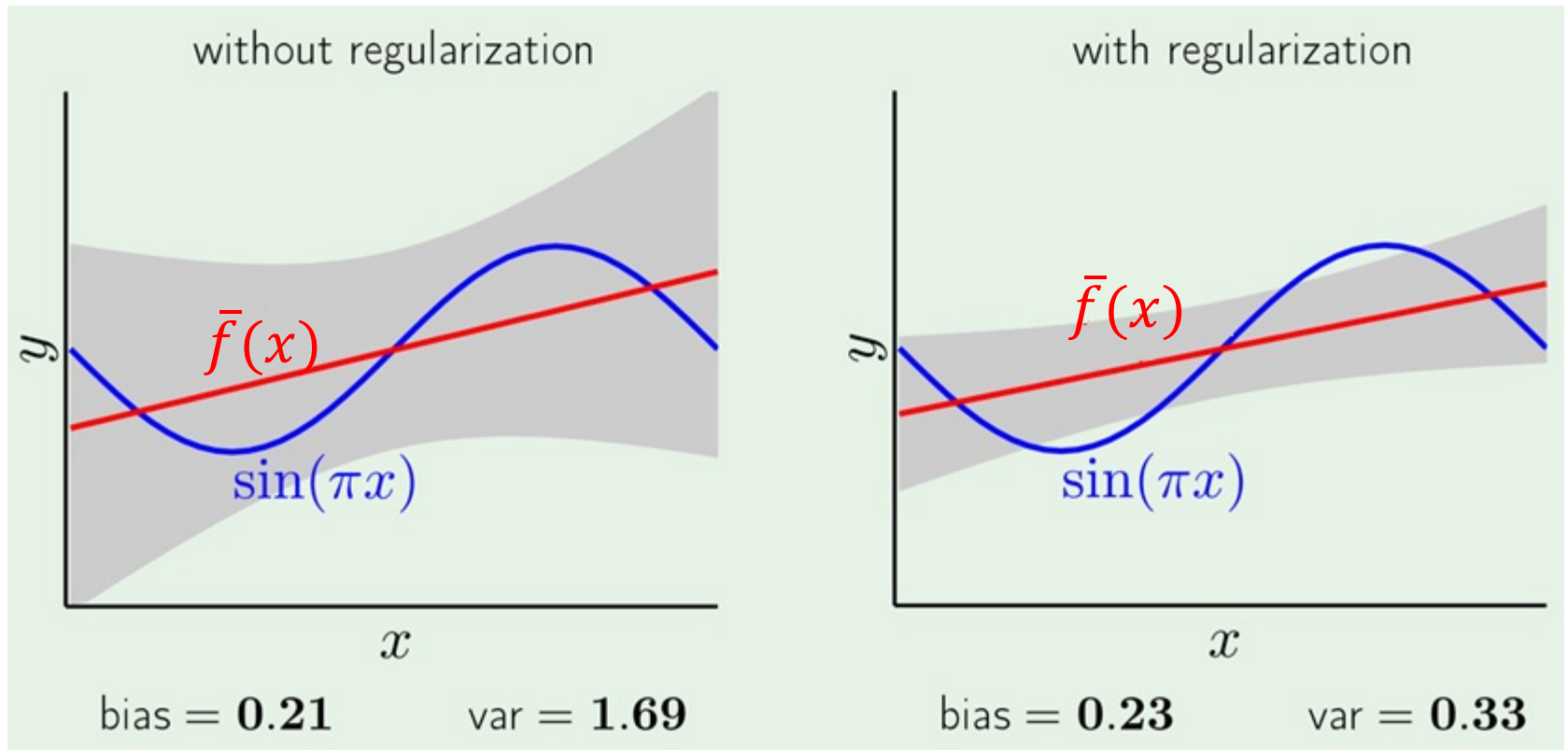


without regularization

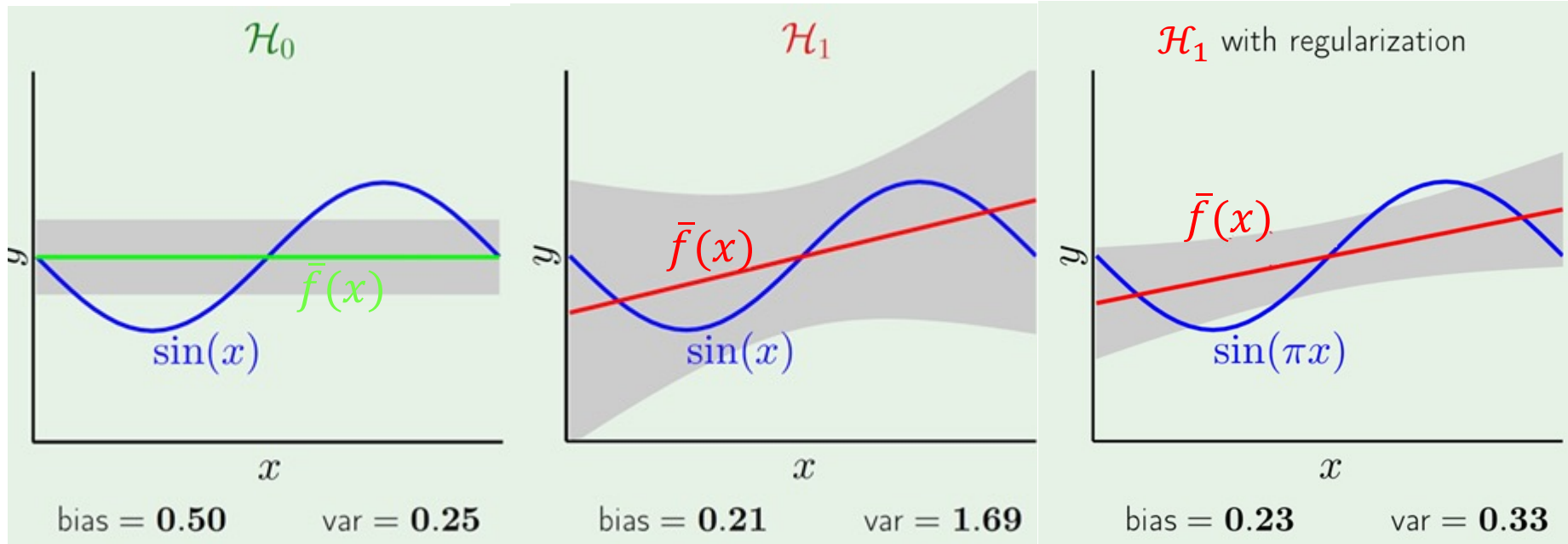


with regularization

# Example II: regularization & bias-variance

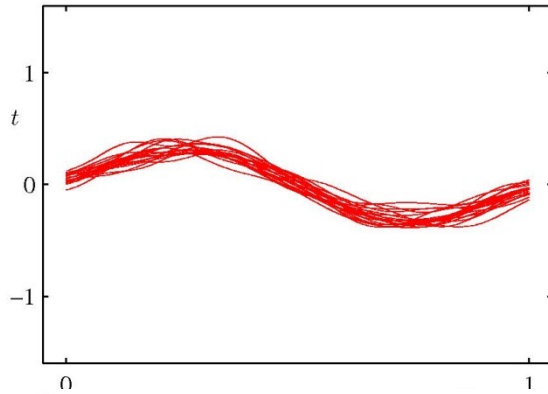


# Winner of $\mathcal{H}_0$ , $\mathcal{H}_1$ , and $\mathcal{H}_1$ with regularization

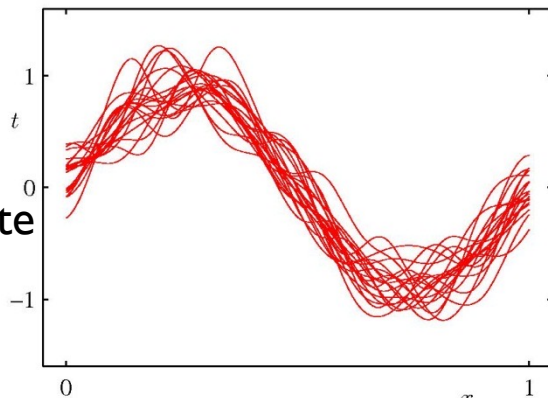


# Example II: regularization & bias/variance

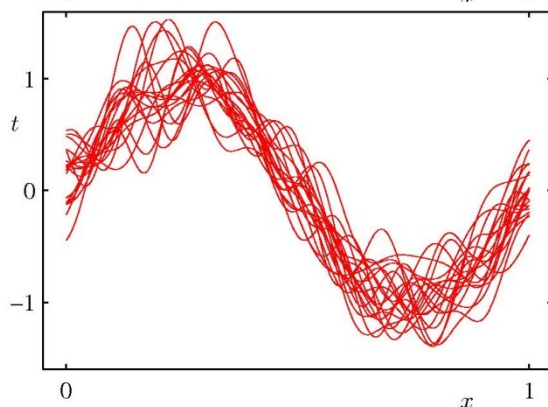
$\lambda$  is large



$\lambda$  is intermediate



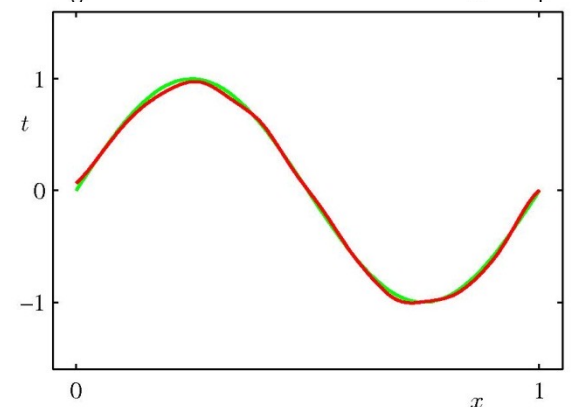
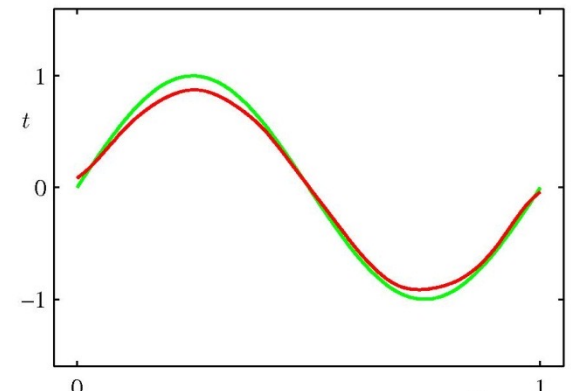
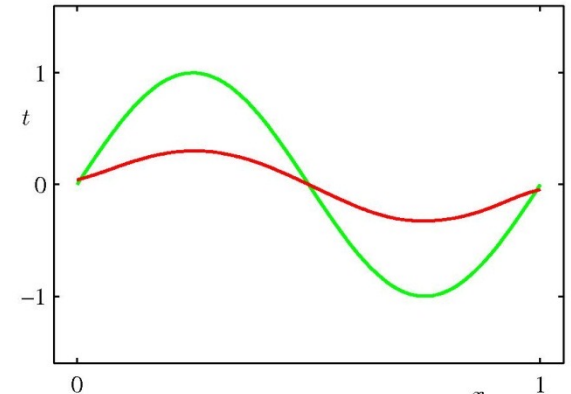
$\lambda$  is small



$L = 100$  data sets

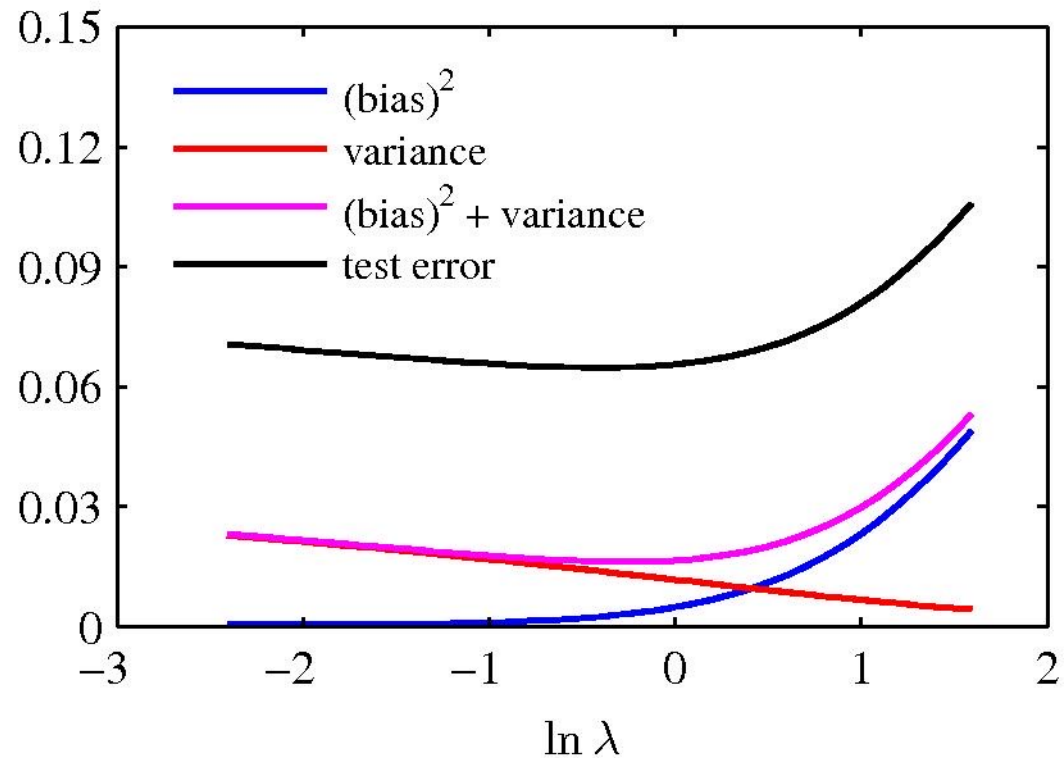
$N = 25$

$m = 25$



[Bishop]

## Example II: Learning curves of bias, variance, and noise



[Bishop]



# Summary

---

- Generalized models
- Overfitting problem & how to avoid it
  - Evaluation and model selection
  - Regularization
- Bias-variance trade-off in regression problem



# 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.

# 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  $\mathbf{x}$ :

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

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

# 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  $\mathbf{x}$ :

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

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

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



# Error decomposition

---

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

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

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

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

---

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

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

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

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

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

# Error decomposition

---

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

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

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

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

$$= \mathbb{E}_{\mathbf{x}} \left[ (f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathbf{x}, y} [(h(\mathbf{x}) - y)^2] \\ + 2 \underbrace{\mathbb{E}_{\mathbf{x}, y} [(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x})) (h(\mathbf{x}) - y)]}_{= 0}$$

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

# Error decomposition

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

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

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

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

$$= \mathbb{E}_{\mathbf{x}} \left[ (f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] + \mathbb{E}_{\mathbf{x}, y} [(h(\mathbf{x}) - y)^2] \\ + 2 \underbrace{\mathbb{E}_{\mathbf{x}, y} [(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x})) (h(\mathbf{x}) - y)]}_{0}$$

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

0



# Error decomposition

---

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

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

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

- Noise shows the irreducible minimum value of the loss function

## Expectation of true error

---

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

# Expectation of true error

---

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

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

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

---

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



## Using the average hypothesis

---

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

## Using the average hypothesis

---

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

# Bias and variance

---

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

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