

Regression and generalization

CE-717: Machine Learning
Sharif University of Technology

M. Soleymani
Fall 2016

Topics

- ▶ Beyond linear regression models
- ▶ Evaluation & model selection
- ▶ Regularization
- ▶ Probabilistic perspective for the regression problem

Recall: Linear regression (squared loss)

- ▶ Linear regression functions

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

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

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

- ▶ Minimizing the squared loss for linear regression

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

- ▶ We obtain $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

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)

Beyond linear regression

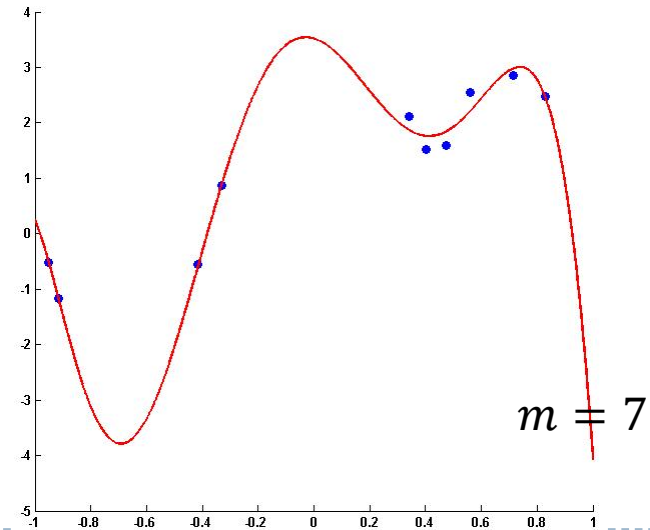
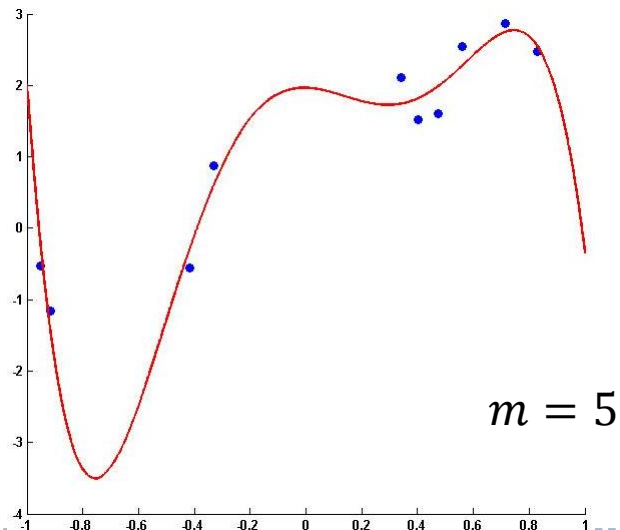
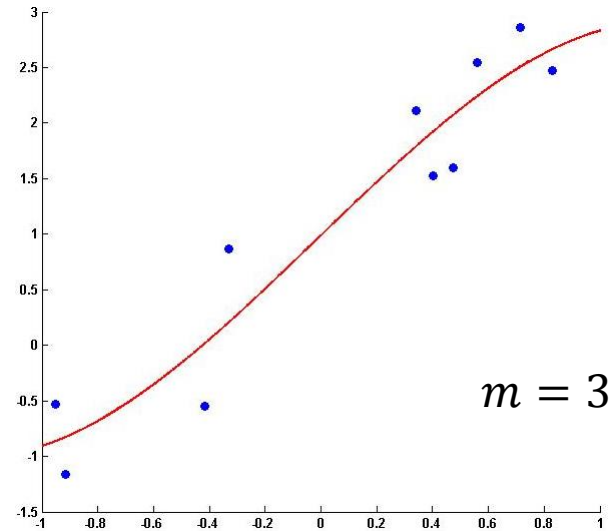
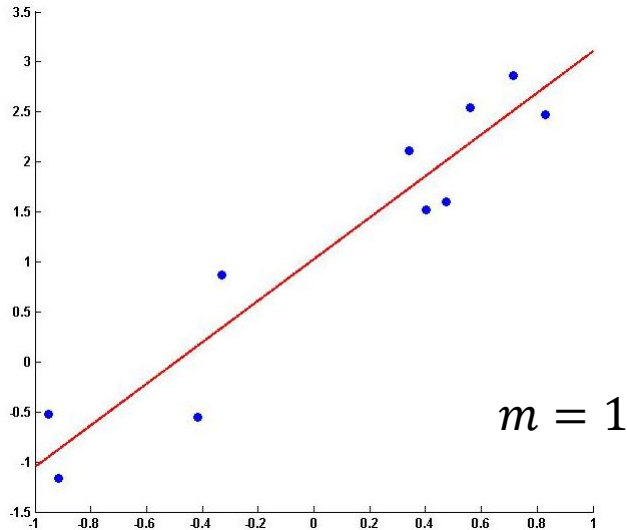
- ▶ m^{th} order polynomial regression (univariate $f : \mathbb{R} \rightarrow \mathbb{R}$)

$$f(x; \mathbf{w}) = w_0 + w_1 x + \dots + w_{m-1} x^{m-1} + w_m x^m$$

- ▶ Solution: $\widehat{\mathbf{w}} = (\mathbf{X}'^T \mathbf{X}')^{-1} \mathbf{X}'^T \mathbf{y}$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X}' = \begin{bmatrix} 1 & x^{(1)1} & x^{(1)2} & \dots & x^{(1)m} \\ 1 & x^{(2)1} & x^{(2)2} & \dots & x^{(2)m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x^{(n)1} & x^{(n)2} & \dots & x^{(n)m} \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} \widehat{w}_0 \\ \widehat{w}_1 \\ \vdots \\ \widehat{w}_m \end{bmatrix}$$

Polynomial regression: example



Generalized linear

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

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

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

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

Basis functions: examples

▶ Linear

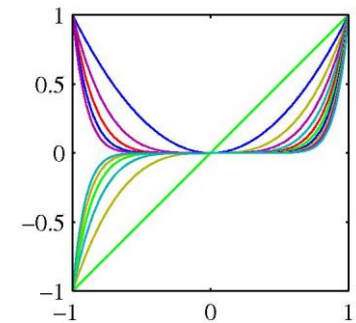
If $m = d$, $\phi_i(\mathbf{x}) = x_i$, $i = 1, \dots, d$, then

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

▶ Polynomial (univariate)

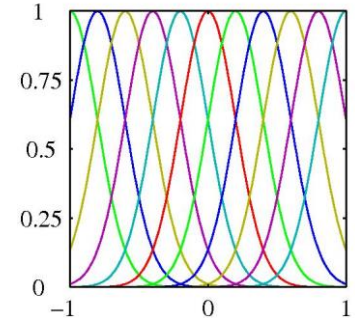
If $\phi_i(x) = x^i$, $i = 1, \dots, m$, then

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

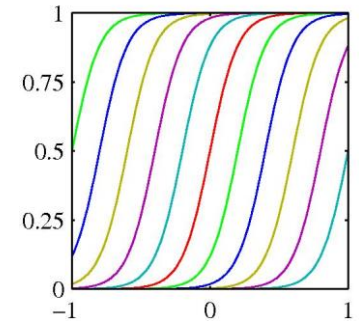


Basis functions: examples

▶ Gaussian: $\phi_j(\mathbf{x}) = \exp\left\{-\frac{(\mathbf{x}-\mathbf{c}_j)^2}{2\sigma_j^2}\right\}$



▶ Sigmoid: $\phi_j(\mathbf{x}) = \sigma\left(\frac{\|\mathbf{x}-\mathbf{c}_j\|}{\sigma_j}\right)$ $\sigma(a) = \frac{1}{1+\exp(-a)}$



Radial Basis Functions: prototypes

- ▶ Predictions based on similarity to “prototypes”:

$$\phi_j(\mathbf{x}) = \exp\left\{-\frac{1}{2\sigma_j^2}\|\mathbf{x} - \mathbf{c}_j\|^2\right\}$$

- ▶ Measuring the similarity to the prototypes $\mathbf{c}_1, \dots, \mathbf{c}_m$
 - ▶ σ^2 controls how quickly it vanishes as a function of the distance to the prototype.
 - ▶ Training examples themselves could serve as prototypes

Generalized linear: optimization

$$\begin{aligned} J(\mathbf{w}) &= \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2 \\ &= \sum_{i=1}^n \left(y^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)}) \right)^2 \end{aligned}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \boldsymbol{\Phi} = \begin{bmatrix} 1 & \phi_1(\mathbf{x}^{(1)}) & \cdots & \phi_m(\mathbf{x}^{(1)}) \\ 1 & \phi_1(\mathbf{x}^{(2)}) & \cdots & \phi_m(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}^{(n)}) & \cdots & \phi_m(\mathbf{x}^{(n)}) \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$$

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

Model complexity and overfitting

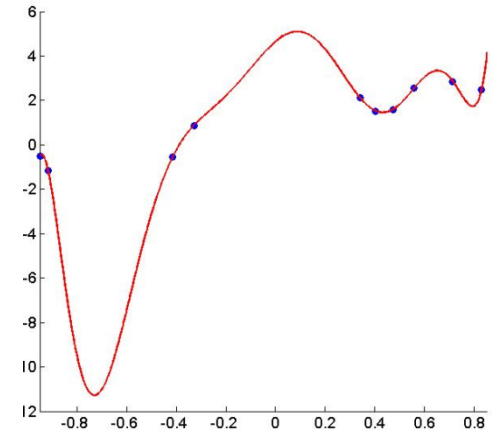
- ▶ With limited training data, models may achieve zero training error but a large test error.

Training
(empirical) loss

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

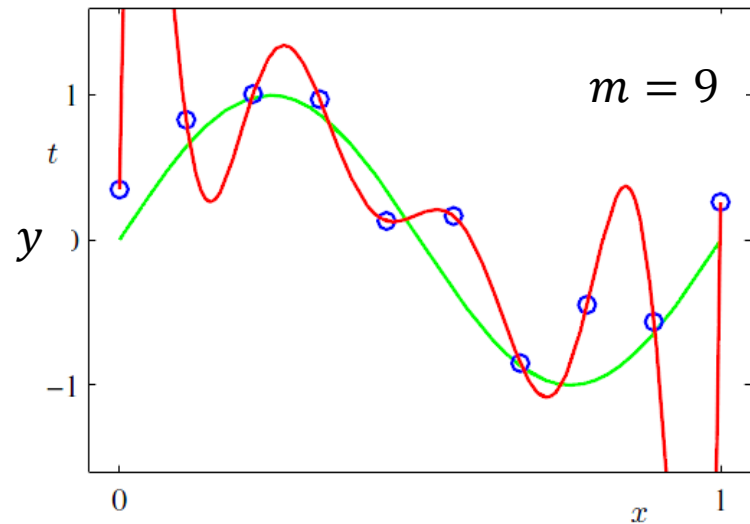
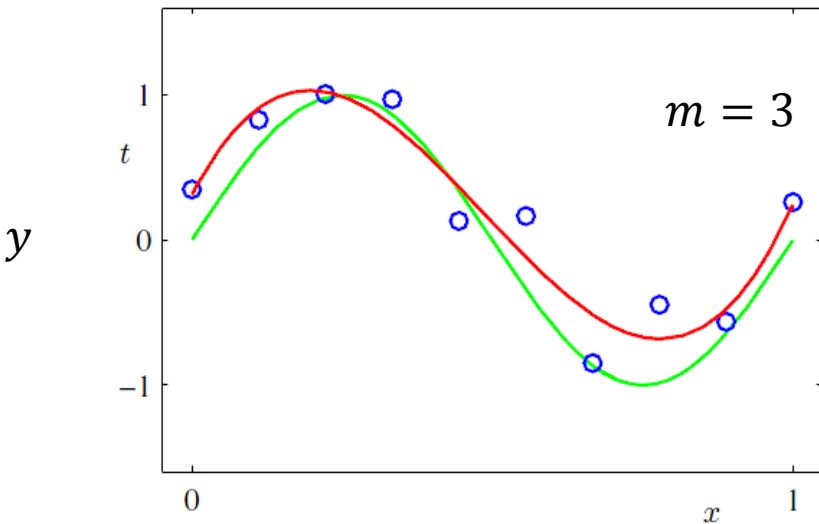
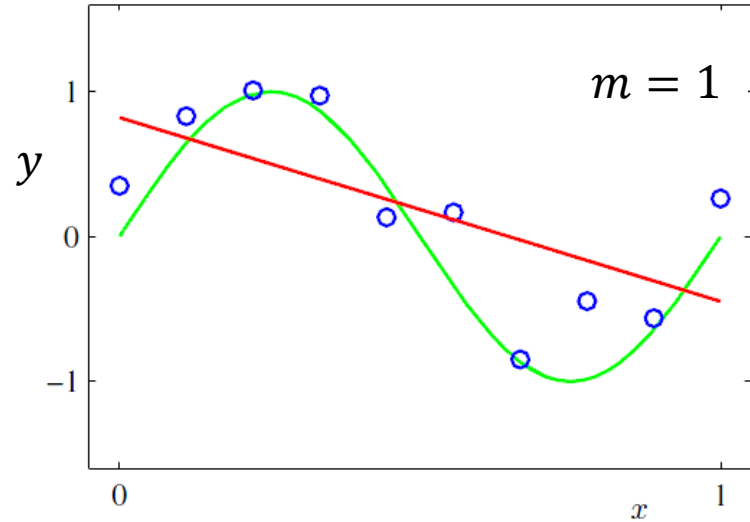
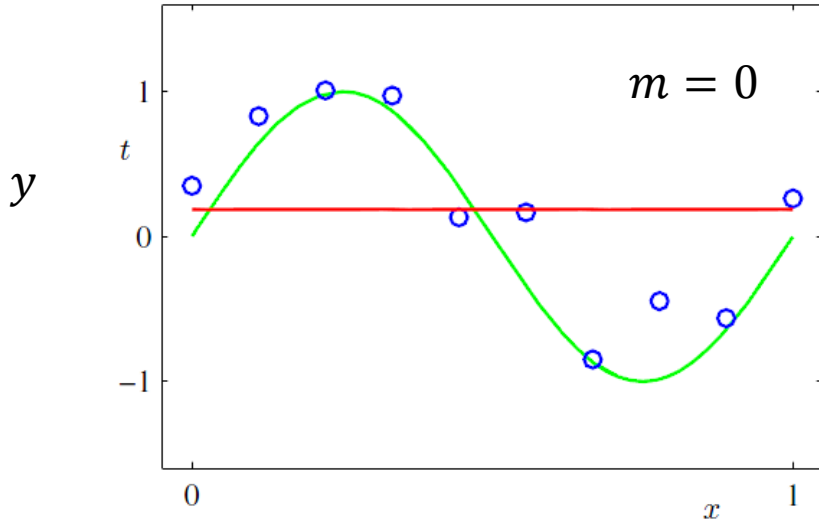
Expected
(test) loss

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



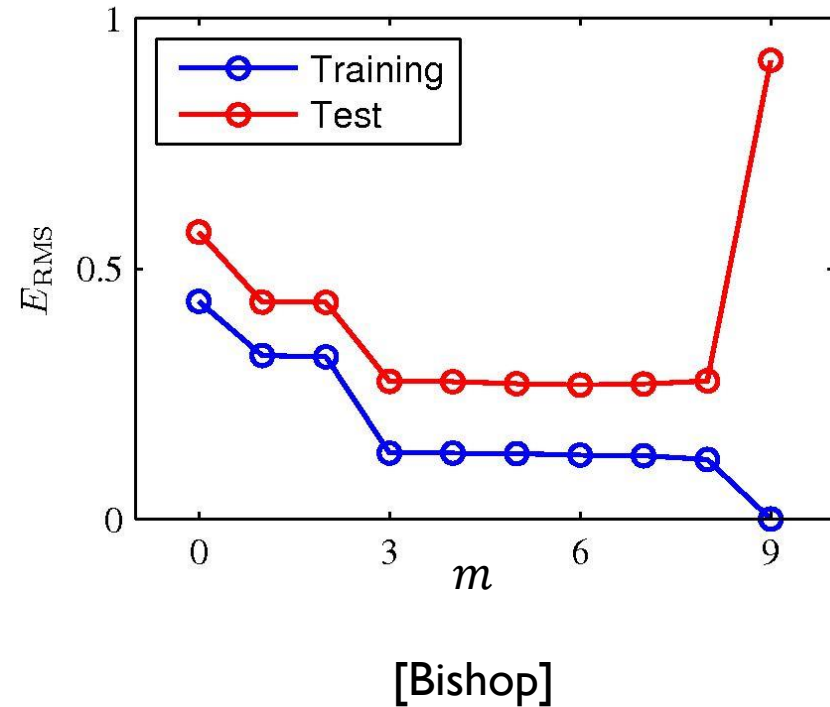
- ▶ Over-fitting: when the training loss no longer bears any relation to the test (generalization) loss.
 - ▶ Fails to generalize to unseen examples.

Polynomial regression



Polynomial regression: training and test error

$$RMSE = \sqrt{\frac{\sum_{i=1}^n \left(y^{(i)} - f(x^{(i)}; \theta) \right)^2}{n}}$$



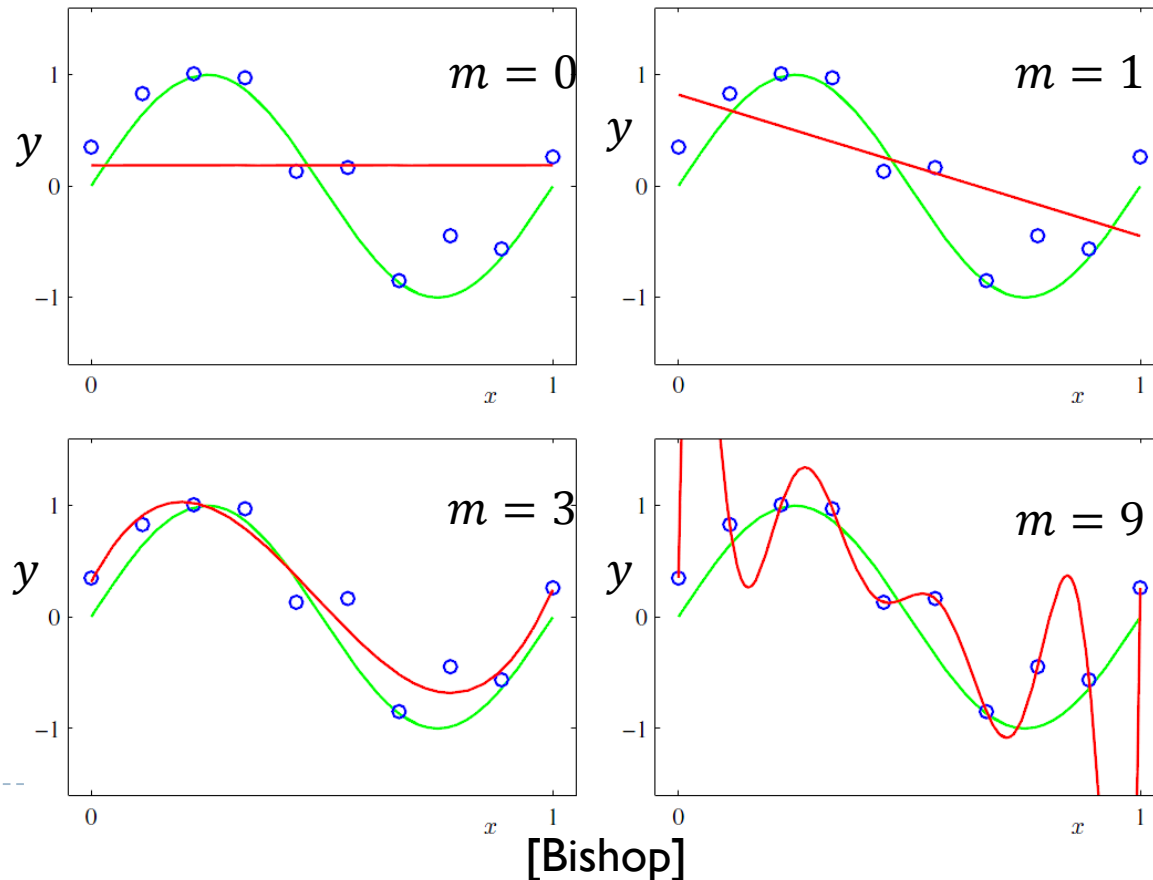
Over-fitting causes

- ▶ **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

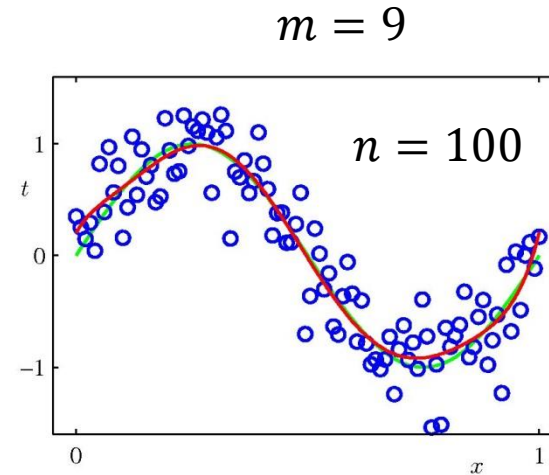
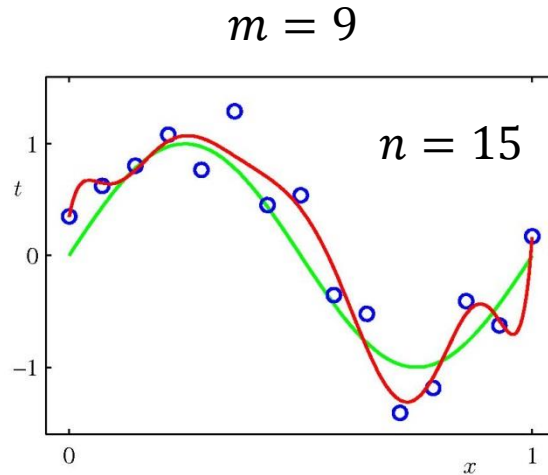
▶ Example:

- ▶ 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]

How to evaluate the learner's performance?

- ▶ Generalization error: true (or expected) error that we would like to optimize
- ▶ Two ways to assess the generalization error is:
 - ▶ Practical: Use a separate data set to test the model
 - ▶ Theoretical: Law of Large numbers
 - ▶ statistical bounds on the difference between training and expected errors

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

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
- ▶ Bayesian approach

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

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

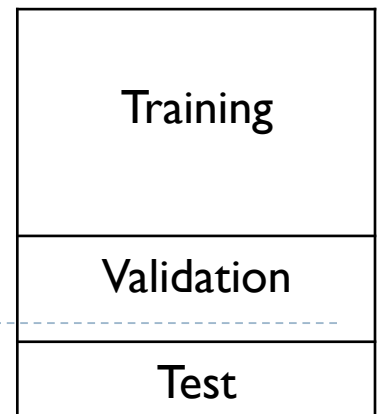
- ▶ Choose the best model based on the validation set error

▶ 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.

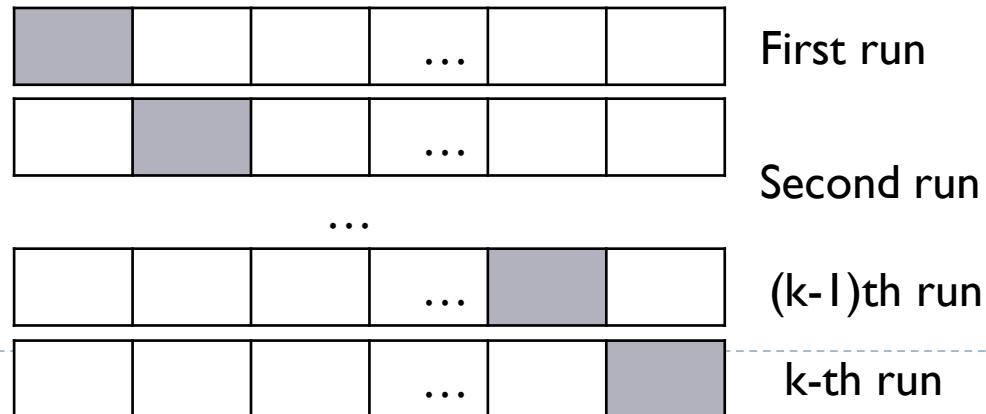
Simple hold out: training, validation, and test sets

- ▶ Simple hold-out chooses the model that minimizes error on validation set.
- ▶ $J_v(\hat{\mathbf{W}})$ 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
 - ▶ Performance scores of the model from k runs are **averaged**.
 - ▶ The average error rate can be considered as an estimation of the true performance.

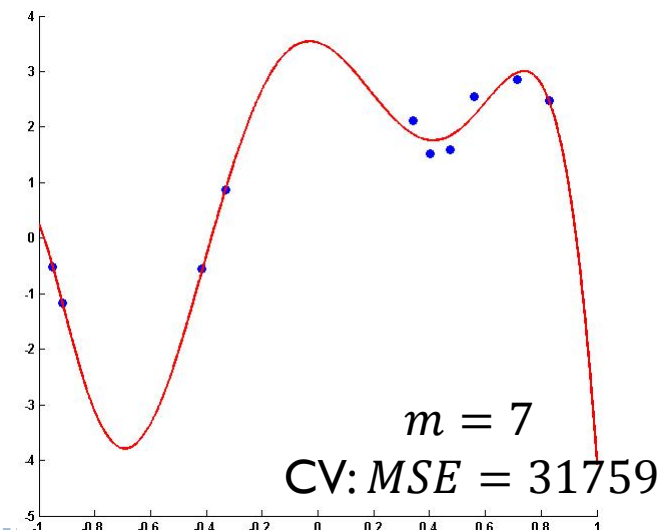
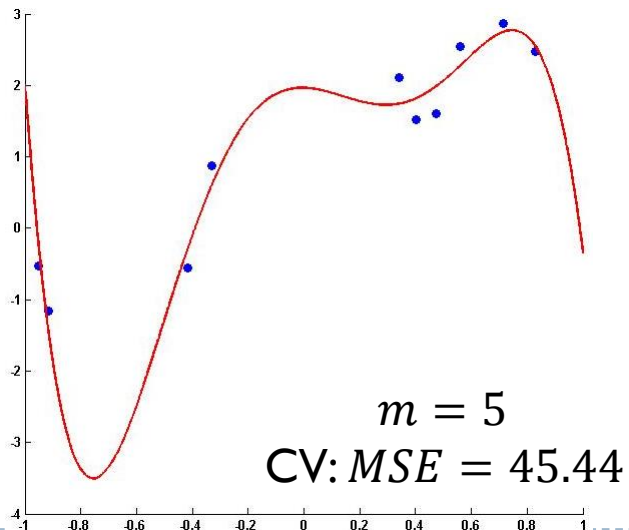
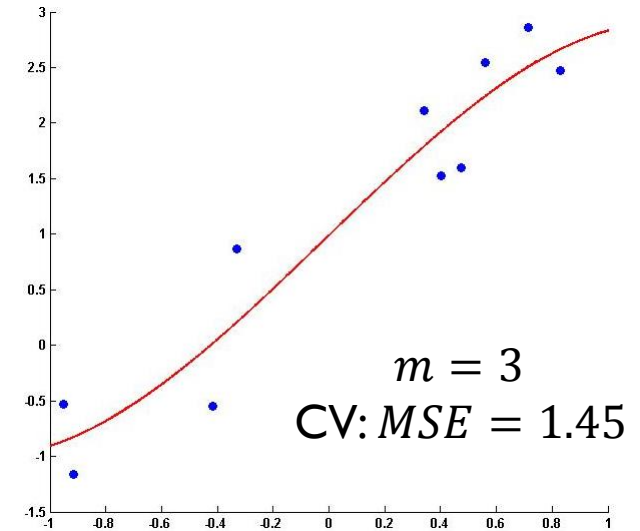
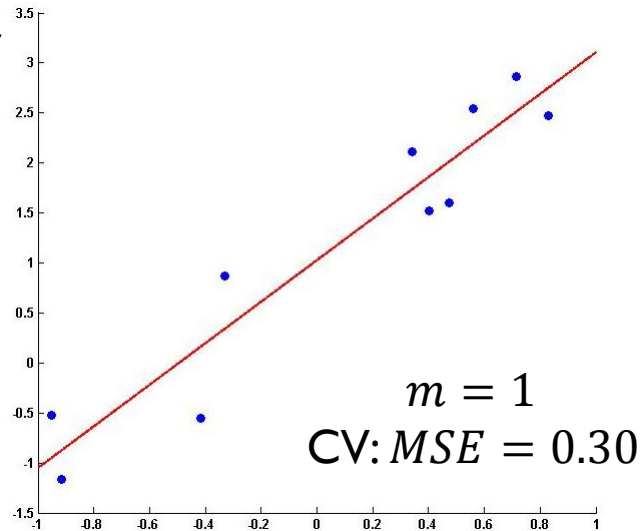


Cross-Validation (CV): Model Selection

- ▶ For each model we first find the average error found 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):

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

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

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 .

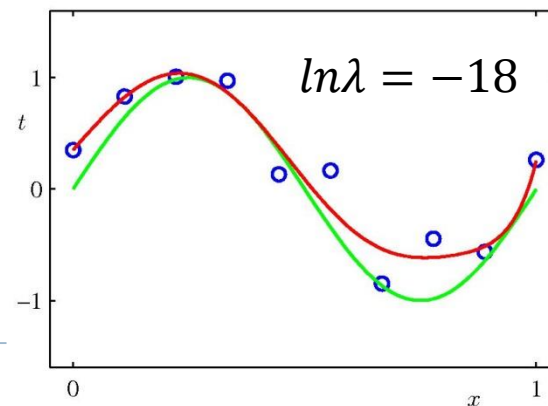
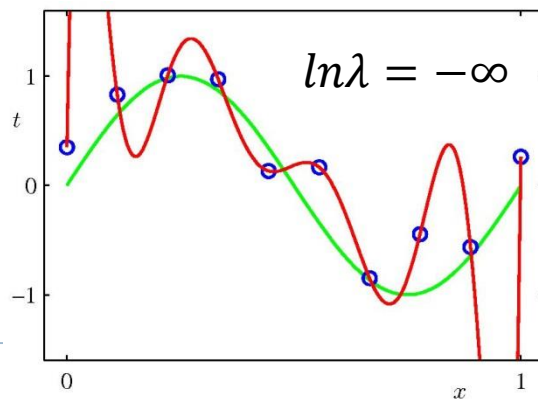
	$M = 0$	$M = 1$	$M = 6$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

[Bishop]

Regularization parameter

	$m = 9$		
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
\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

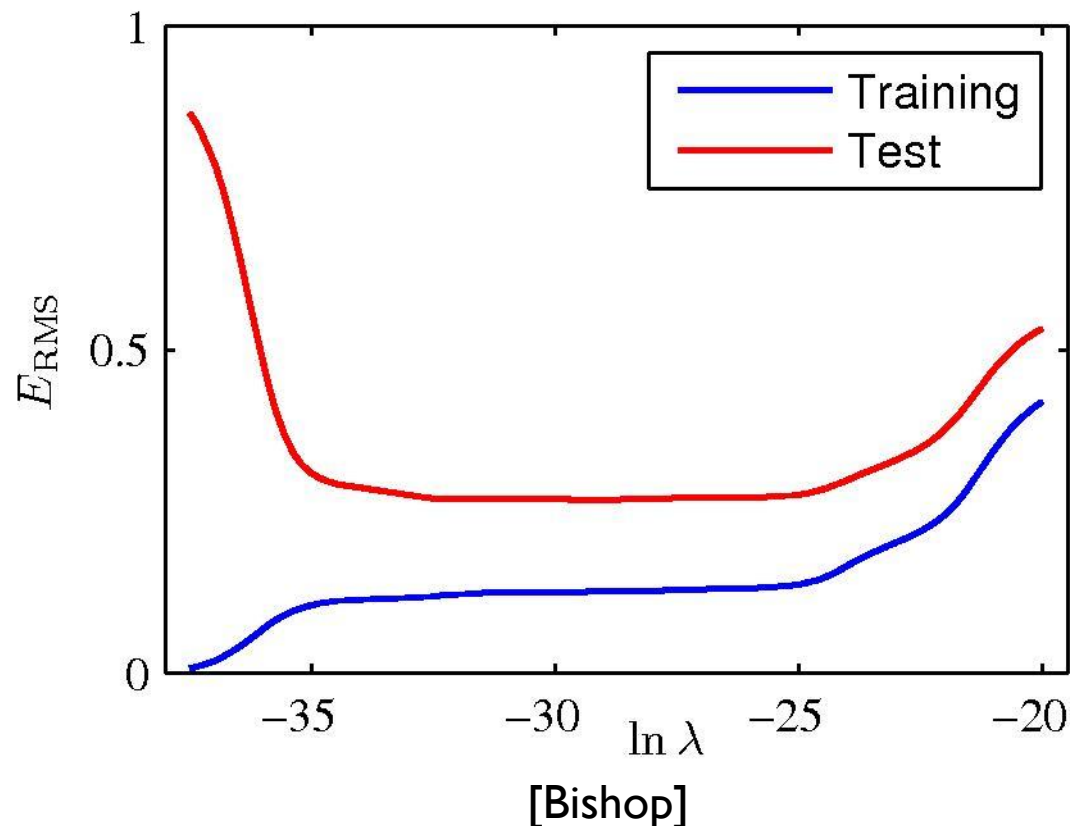
[Bishop]



Regularization parameter

► Generalization

- λ 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 λ .
- ▶ 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

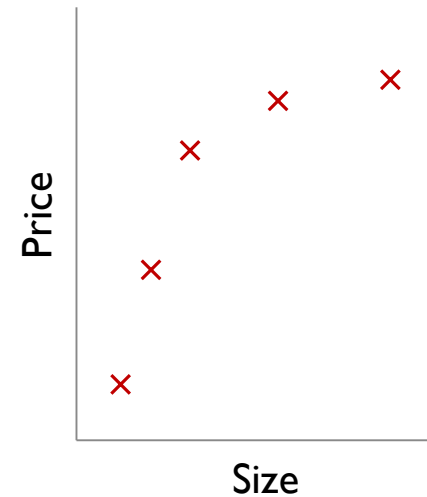


$$w_0 + w_1x$$

Less complex \mathcal{H}



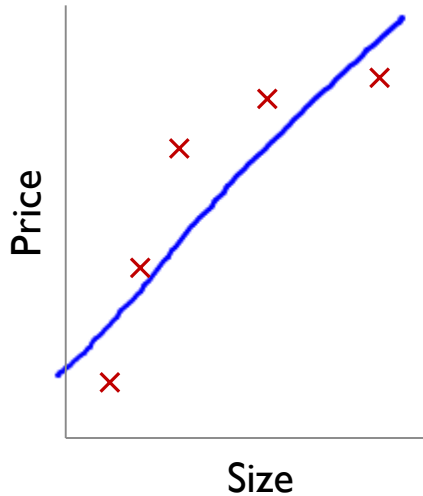
$$w_0 + w_1x + w_2x^2$$



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

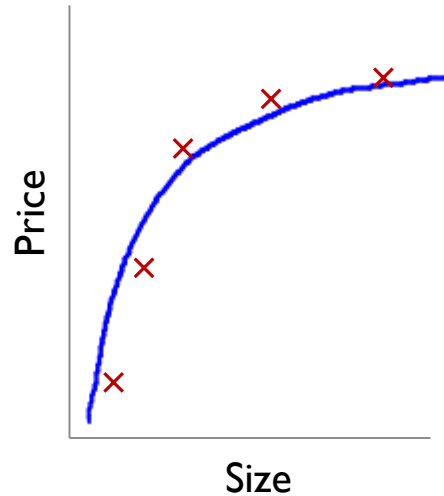
More complex \mathcal{H}

Complexity of Hypothesis Space: Example



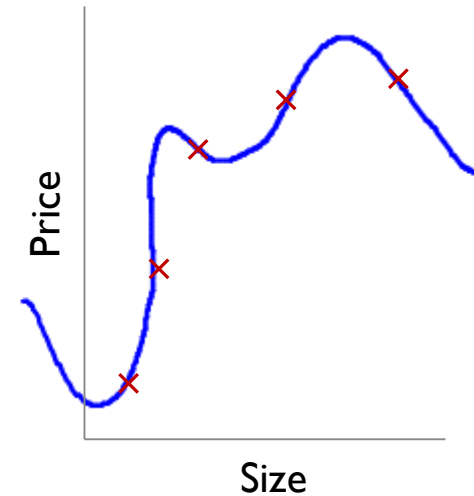
$$w_0 + w_1x$$

Less complex \mathcal{H}



$$w_0 + w_1x + w_2x^2$$

More complex \mathcal{H}

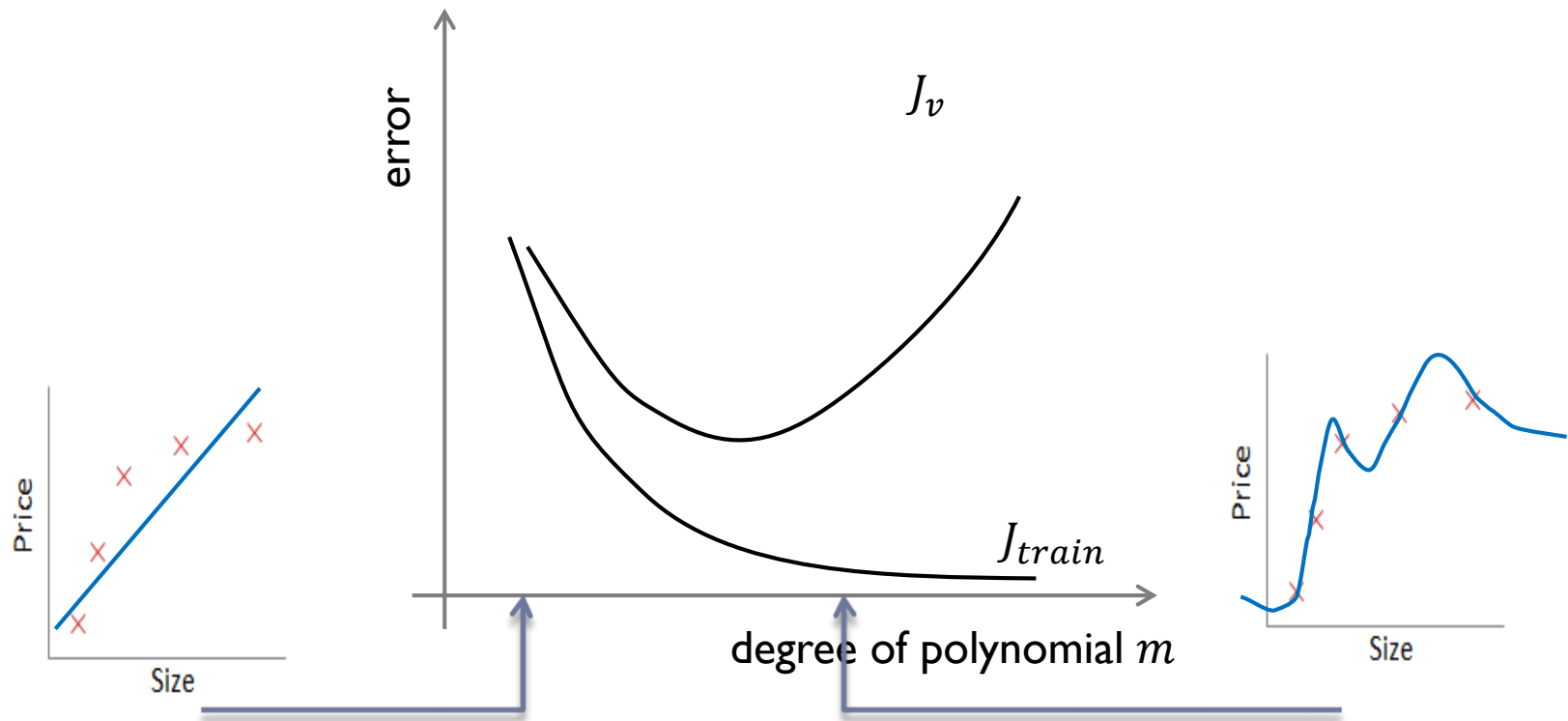


$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

Complexity of Hypothesis Space: Example

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

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



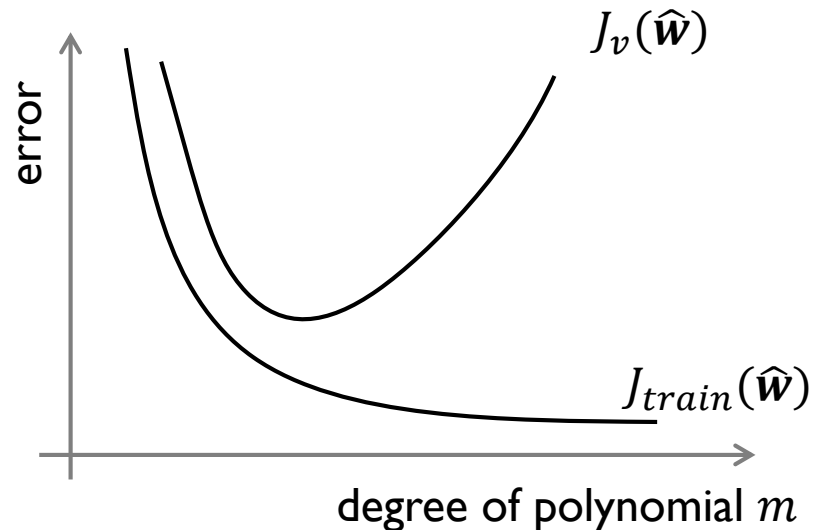
Complexity of Hypothesis Space

▶ 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} :

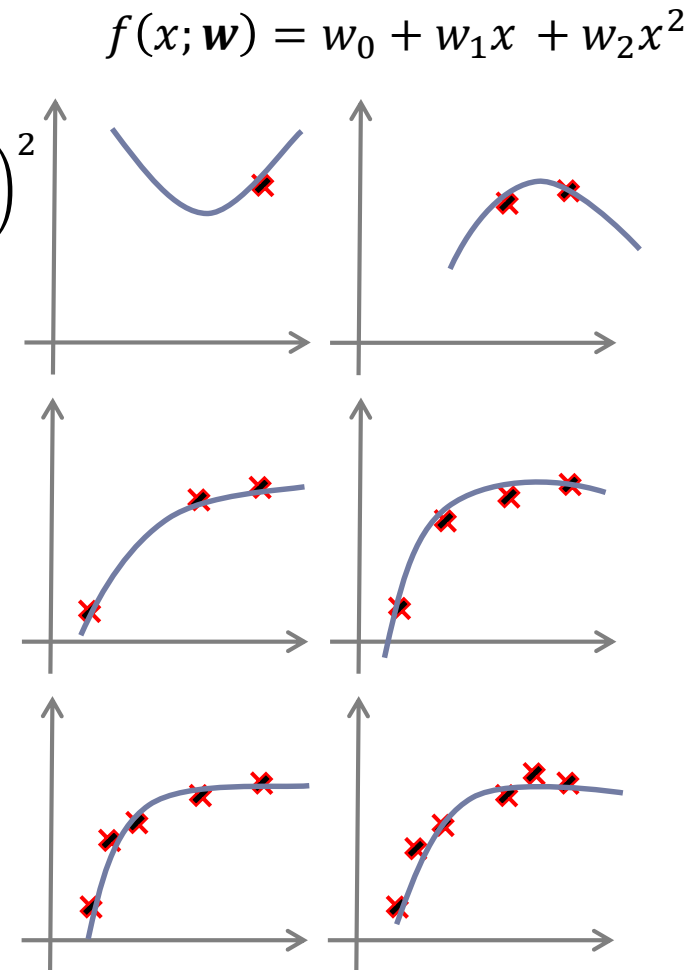
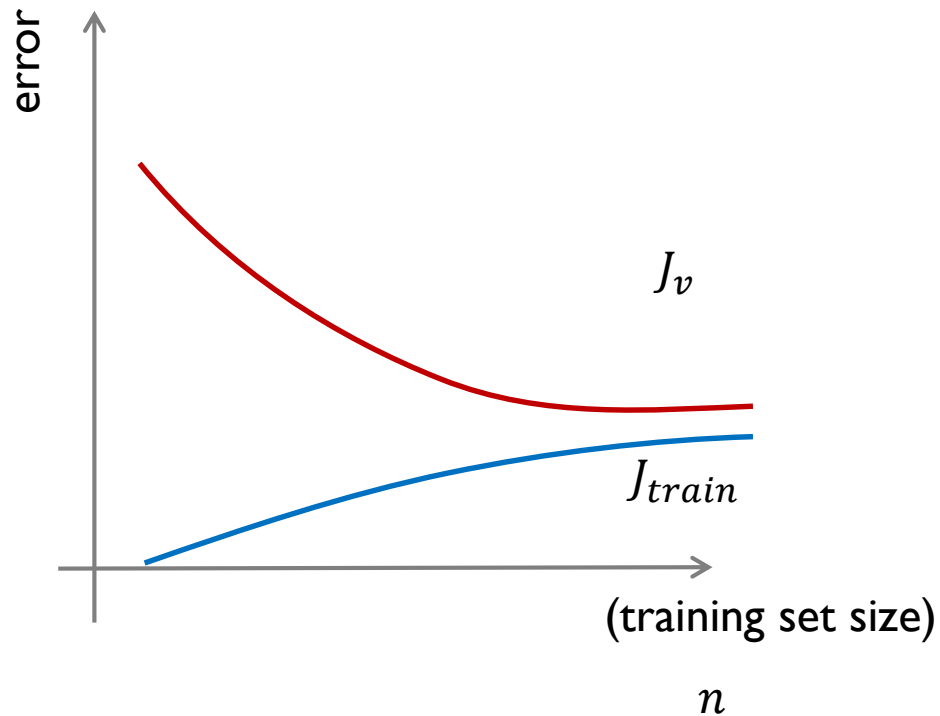
▶ $J_{train}(\hat{\mathbf{w}}) \ll J_v(\hat{\mathbf{w}})$ and $J_{train}(\hat{\mathbf{w}})$ is low



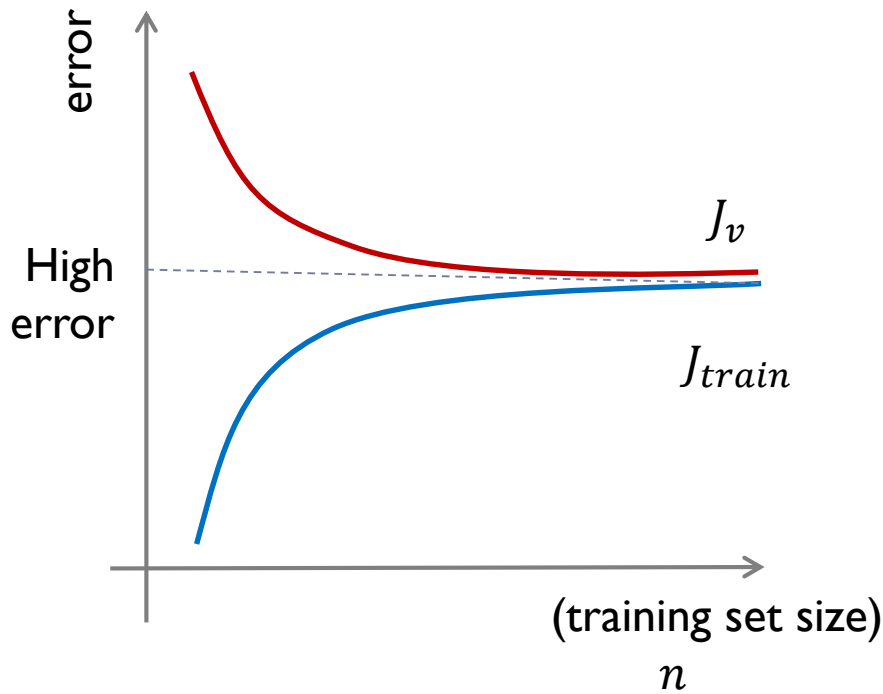
Size of training set

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

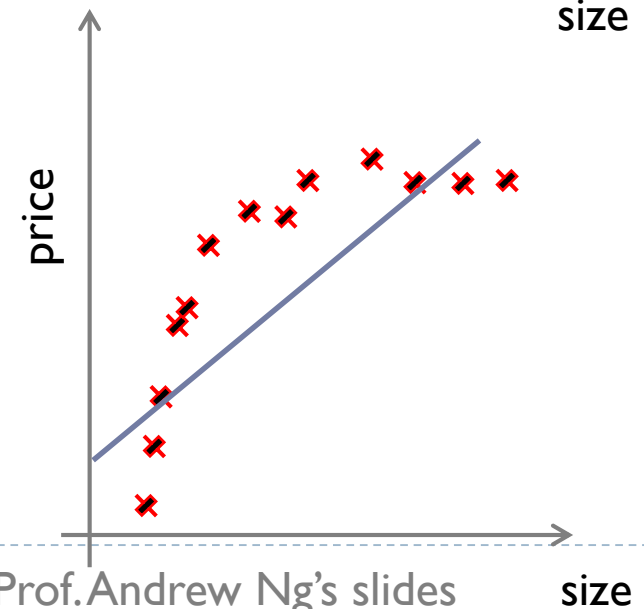
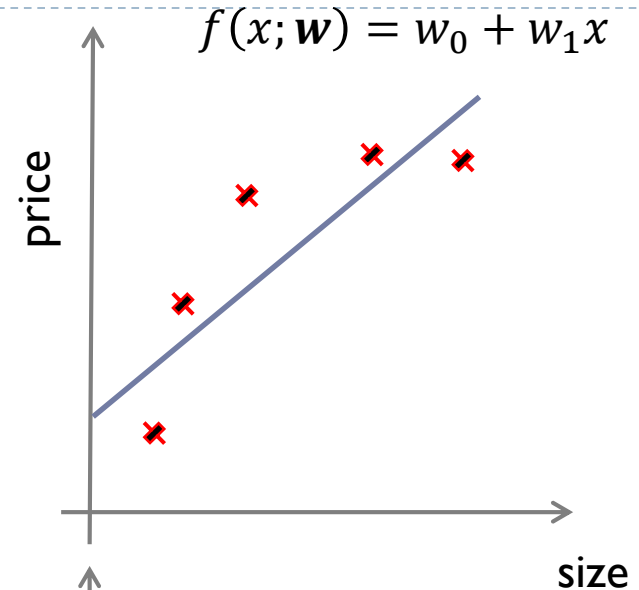
$$J_{\text{train}}(\mathbf{w}) = \frac{1}{n_{\text{train}}} \sum_{i \in \text{train_set}} \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2$$



Less complex \mathcal{H}

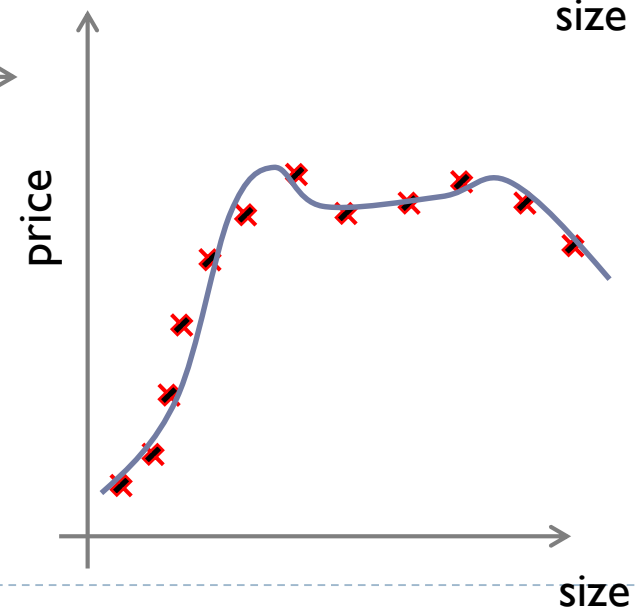
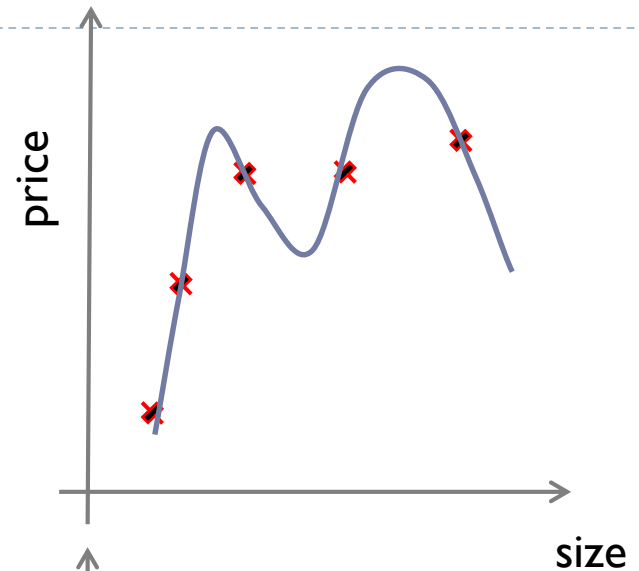
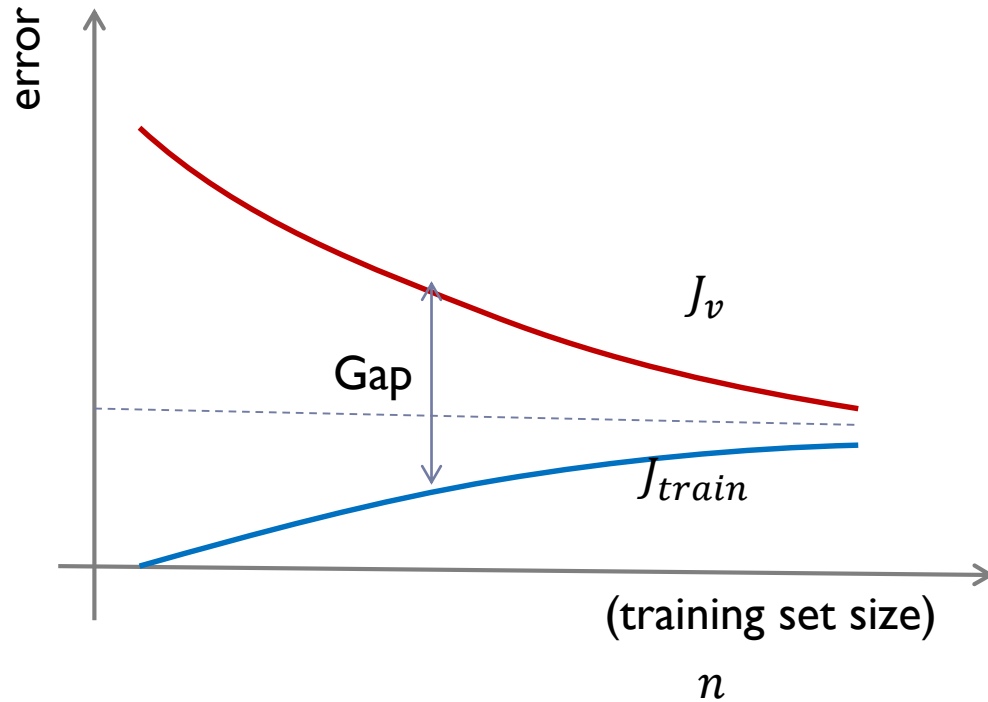


If model is very simple, getting more training data will not (by itself) help much.



More complex \mathcal{H}

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

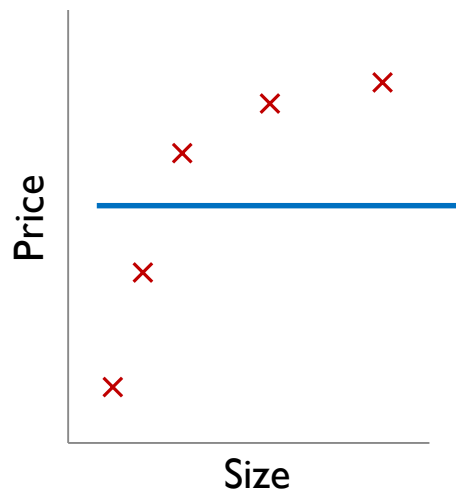


For more complex models, getting more training data is usually helps.

Regularization: Example

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

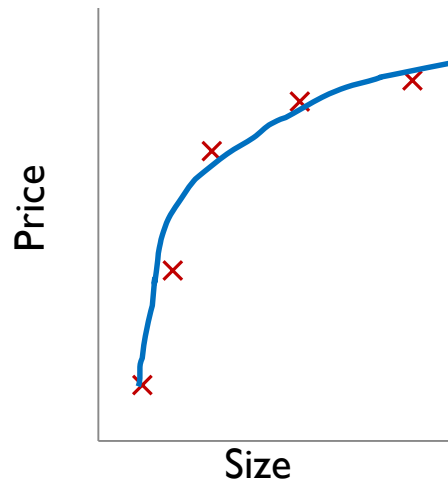
$$J(\mathbf{w}) = \frac{1}{n} \left(\sum_{i=1}^n \left(y^{(i)} - f(x^{(i)}; \mathbf{w}) \right)^2 + \lambda \mathbf{w}^T \mathbf{w} \right)$$



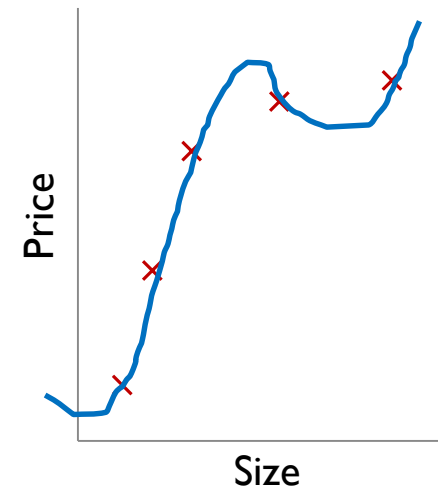
Large λ

(Prefer to more simple models)

$$w_1 = w_2 \approx 0$$



Intermediate λ



Small λ

(Prefer to more complex models)

$$\lambda = 0$$

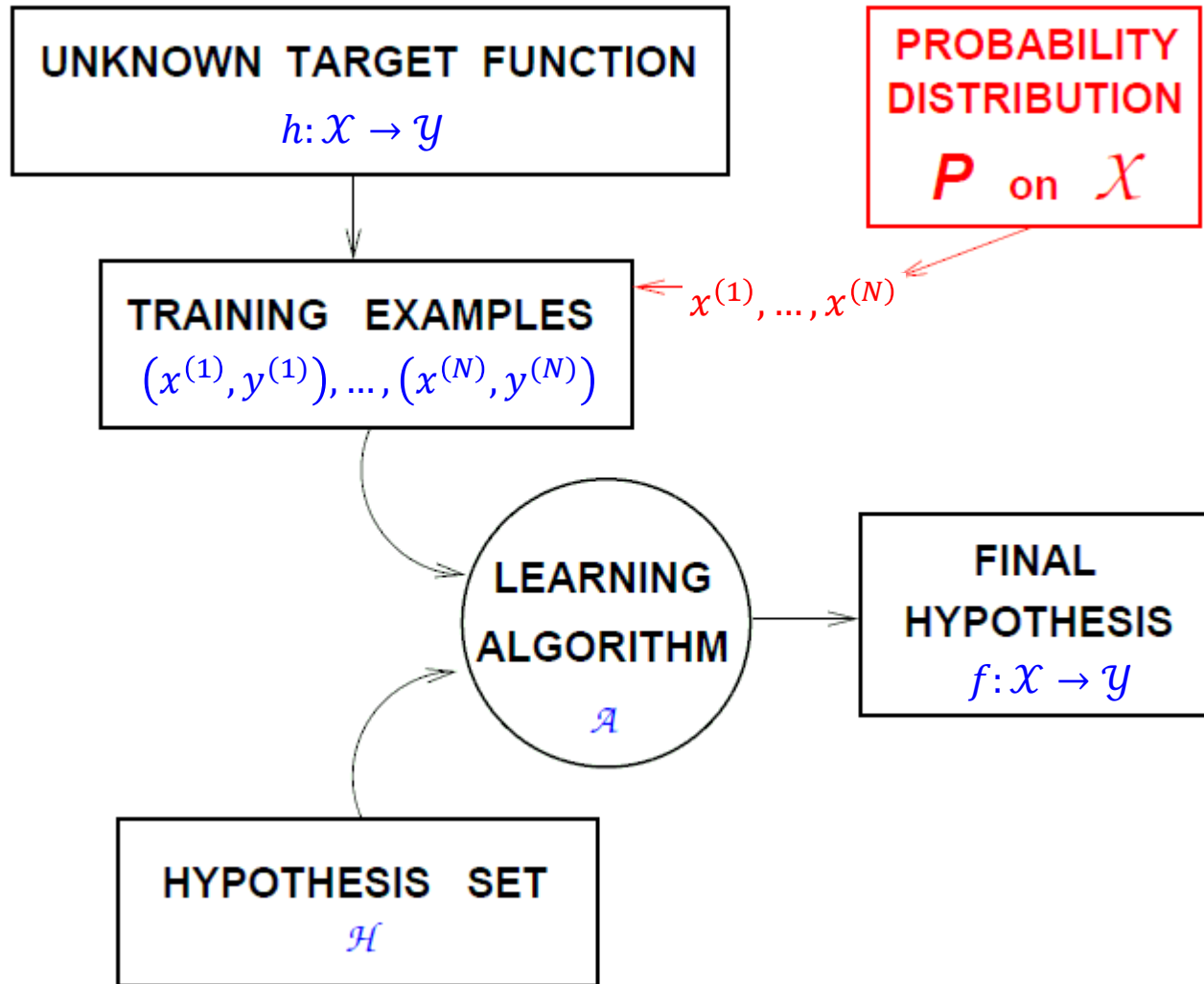
Model complexity: Bias-variance trade-off

- ▶ 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*.

Formal discussion on bias, variance, and noise

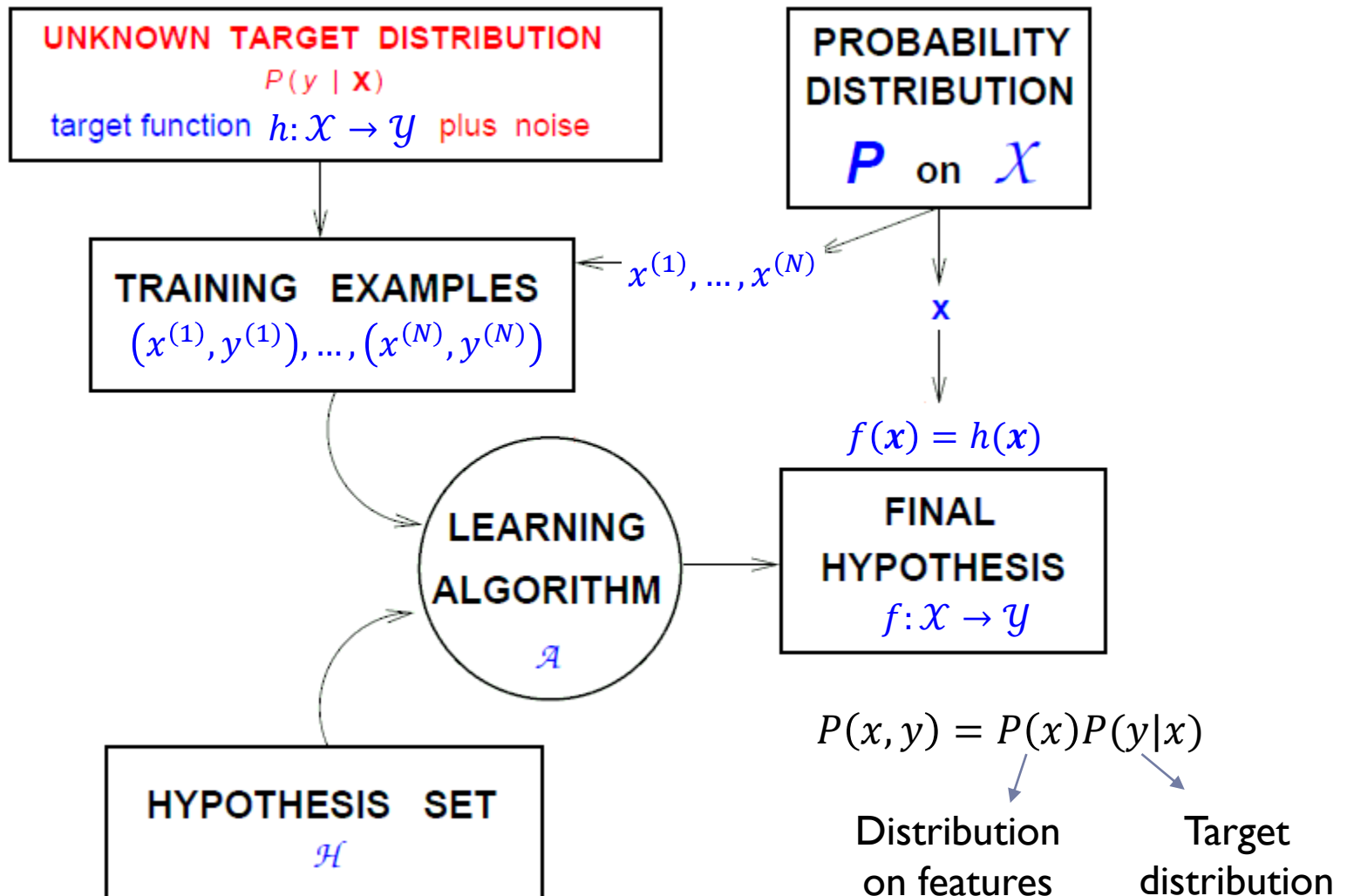
- ▶ Best unrestricted regression function
- ▶ Noise
- ▶ Bias and variance

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} :

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

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

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

Error decomposition

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

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

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

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

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

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

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

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

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

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

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

Bias and variance

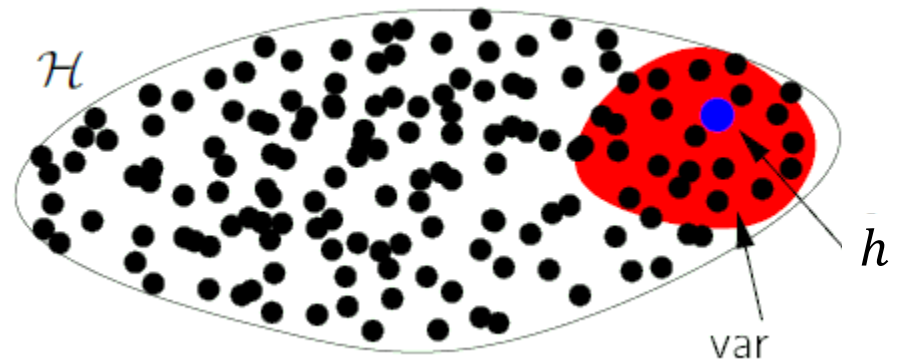
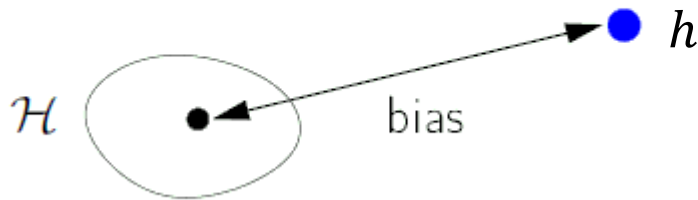
$$\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - h(\mathbf{x}))^2 \right] = \underbrace{\mathbb{E}_{\mathcal{D}} \left[(f_{\mathcal{D}}(\mathbf{x}) - \bar{f}(\mathbf{x}))^2 \right]}_{\text{var}(\mathbf{x})} + \underbrace{(\bar{f}(\mathbf{x}) - h(\mathbf{x}))^2}_{\text{bias}(\mathbf{x})}$$

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

Bias-variance trade-off

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

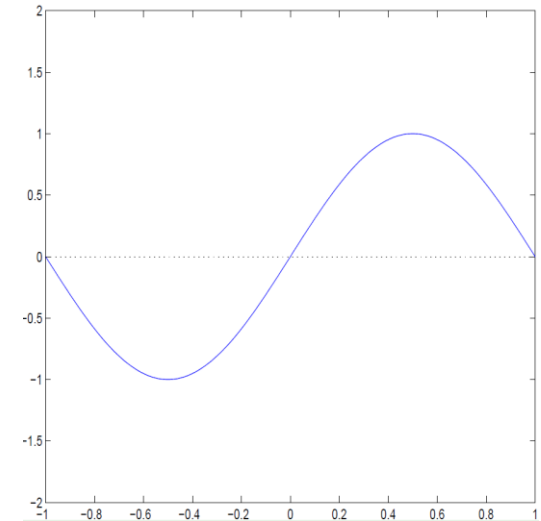
$$\text{bias} = \mathbb{E}_{\mathbf{x}} \left[\bar{f}(\mathbf{x}) - h(\mathbf{x}) \right]$$



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

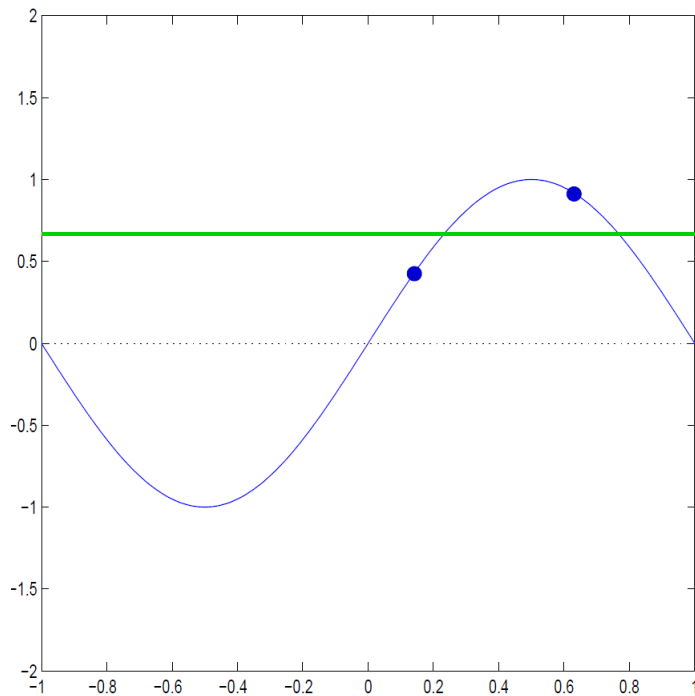
Example: sin target

- ▶ 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 ?

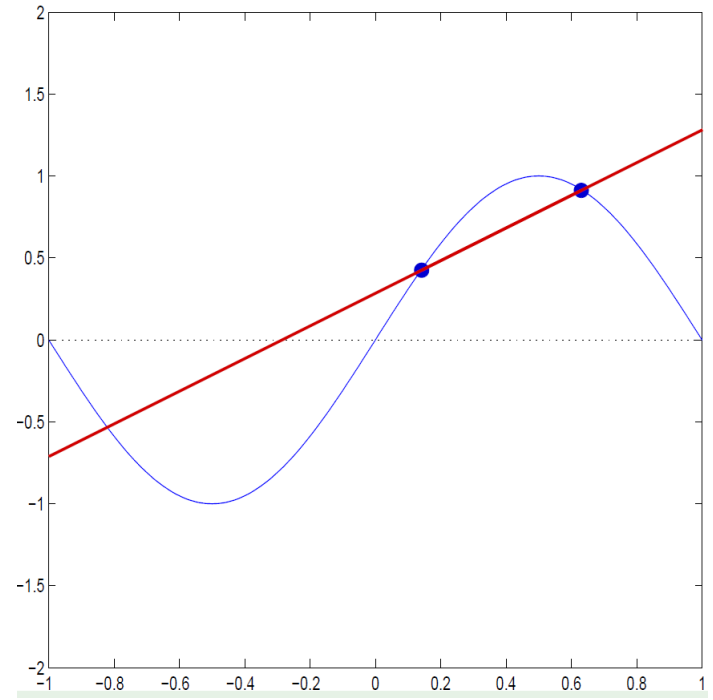


Learning from a training set

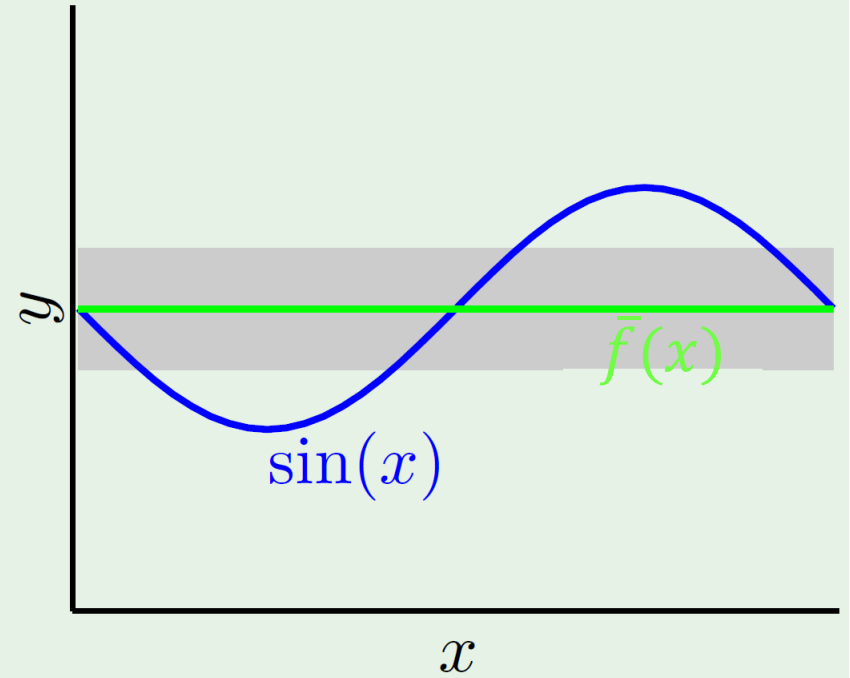
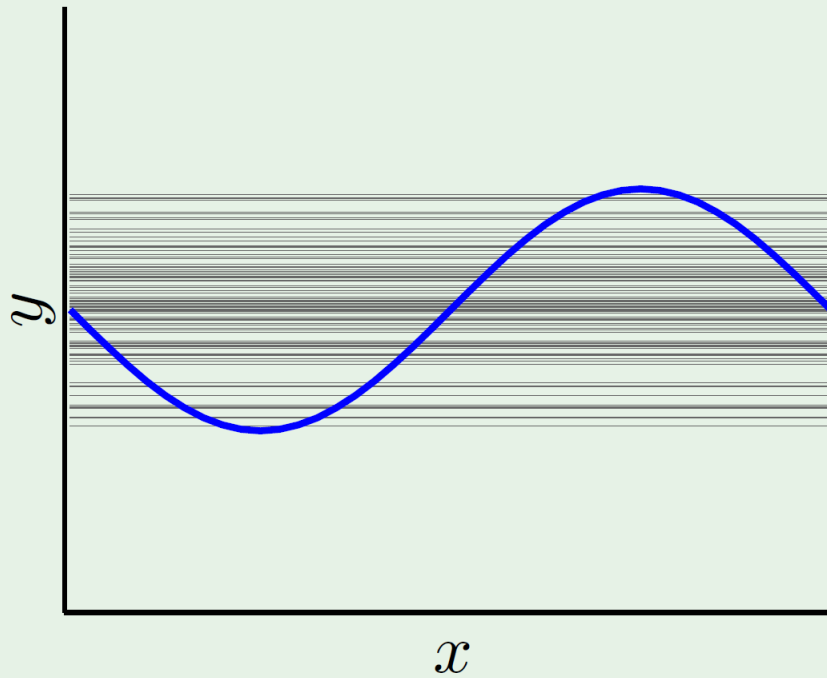
\mathcal{H}_0



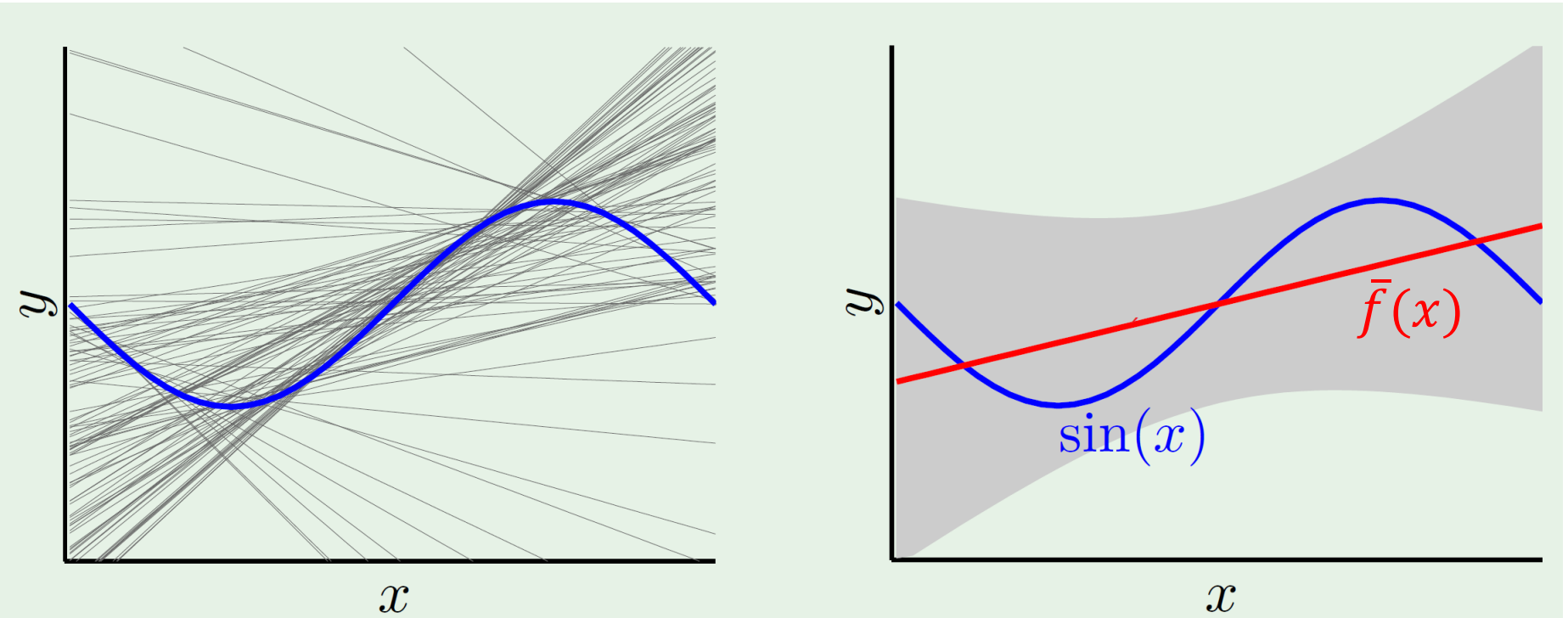
\mathcal{H}_1



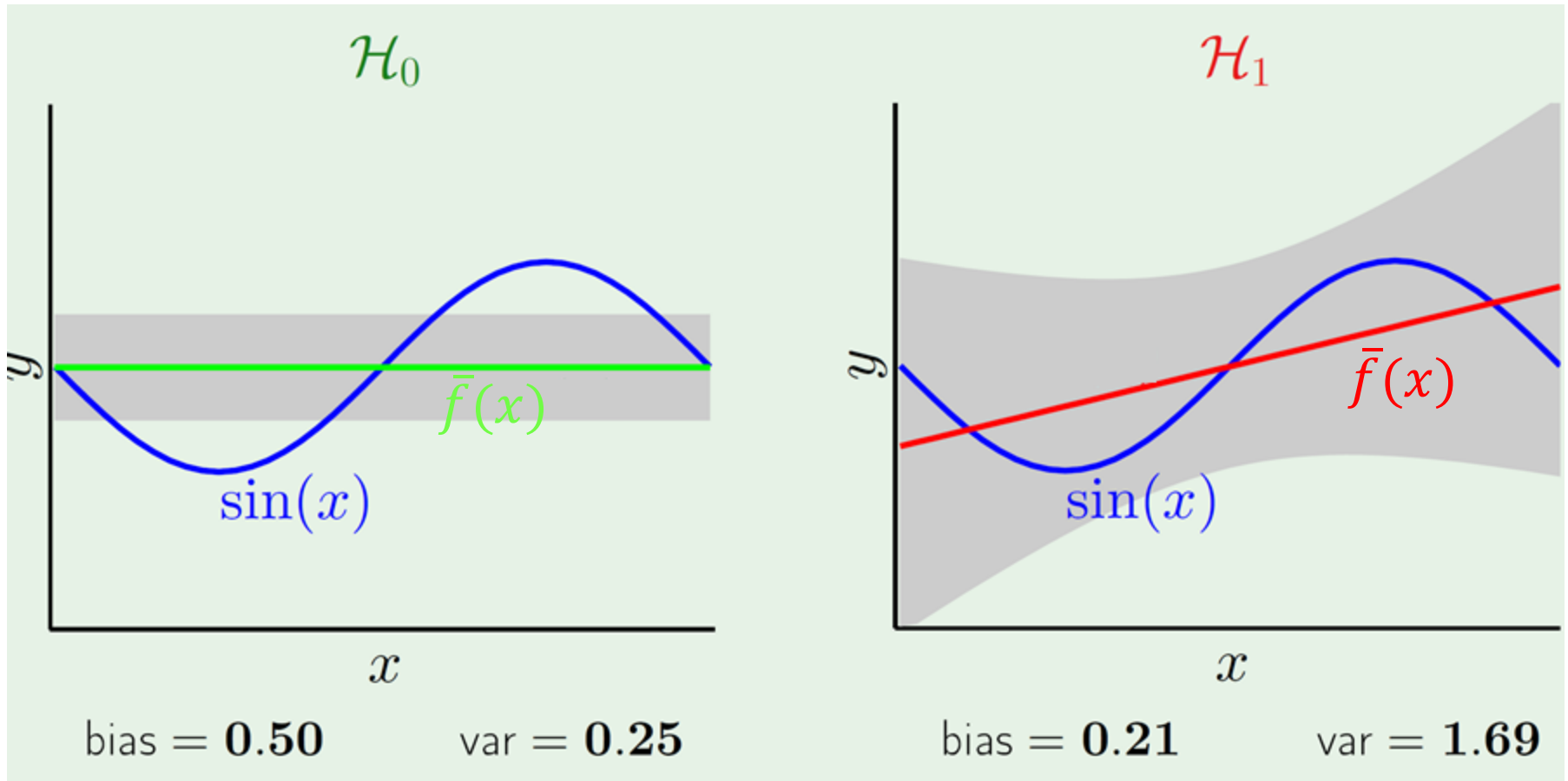
Variance \mathcal{H}_0



Variance \mathcal{H}_1



Which is better?

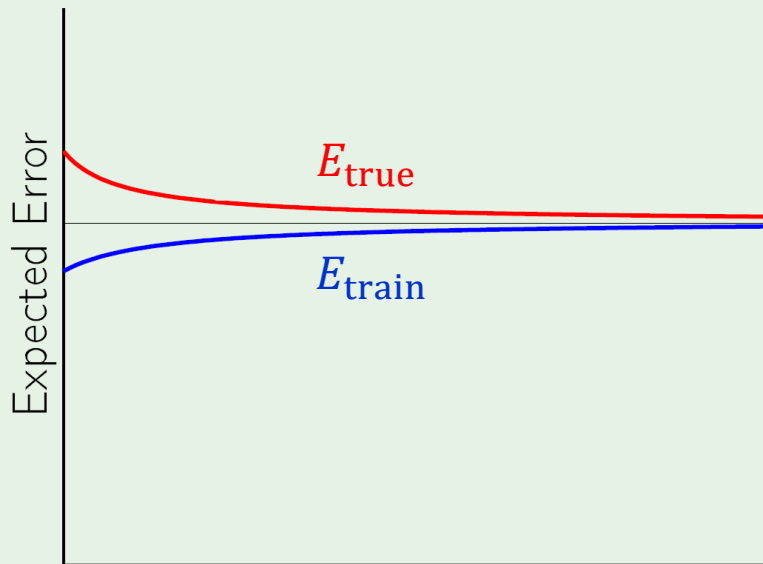


Lesson

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

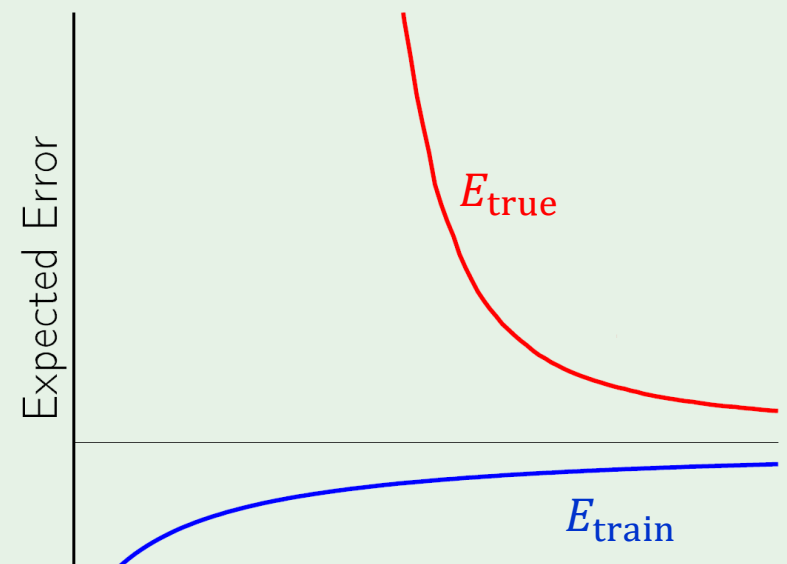
Expected training and true error curves

- ▶ Errors vary with the number of training samples



Number of Data Points, N

Simple Model

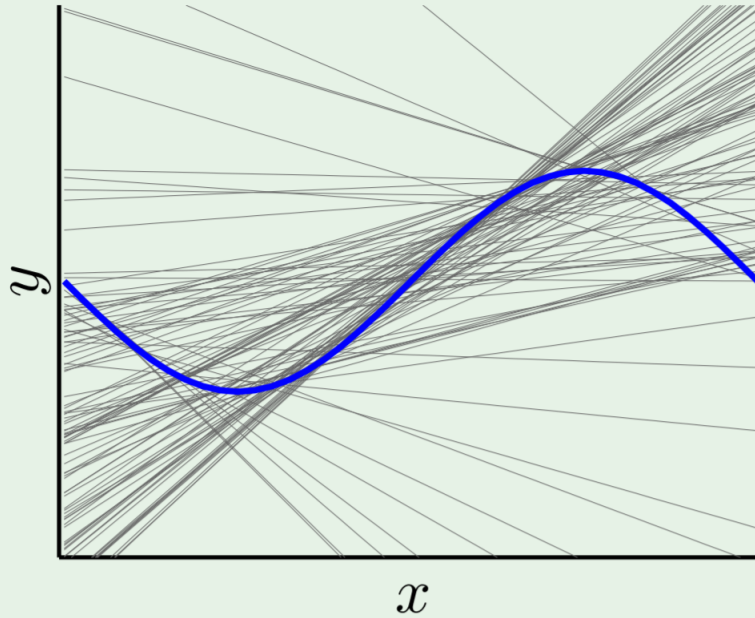


Number of Data Points, N

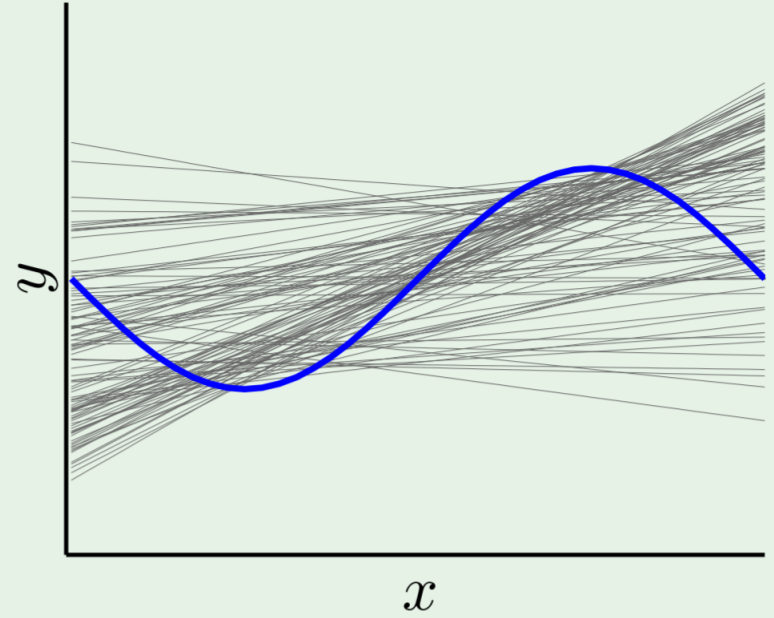
Complex Model

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

Regularization

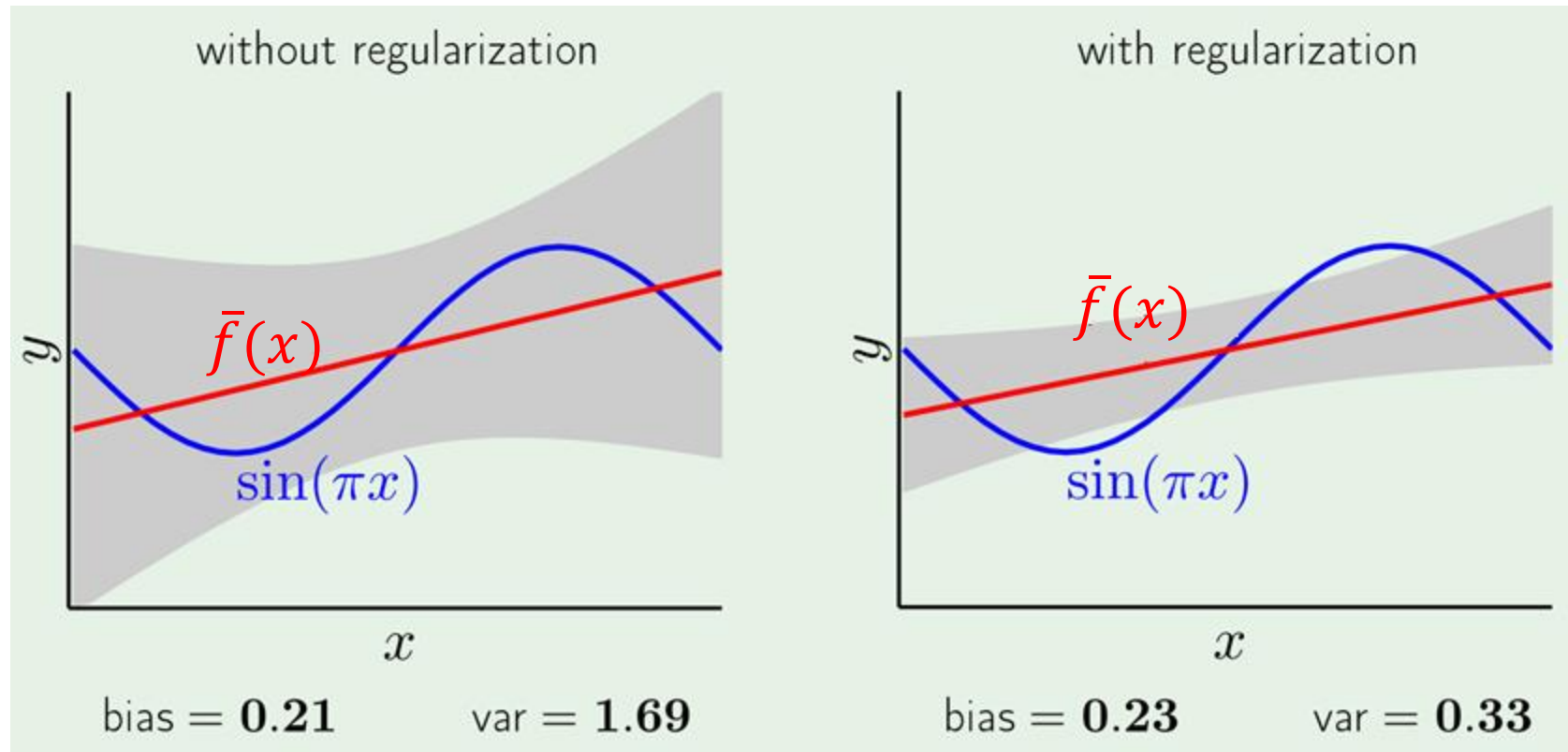


without regularization

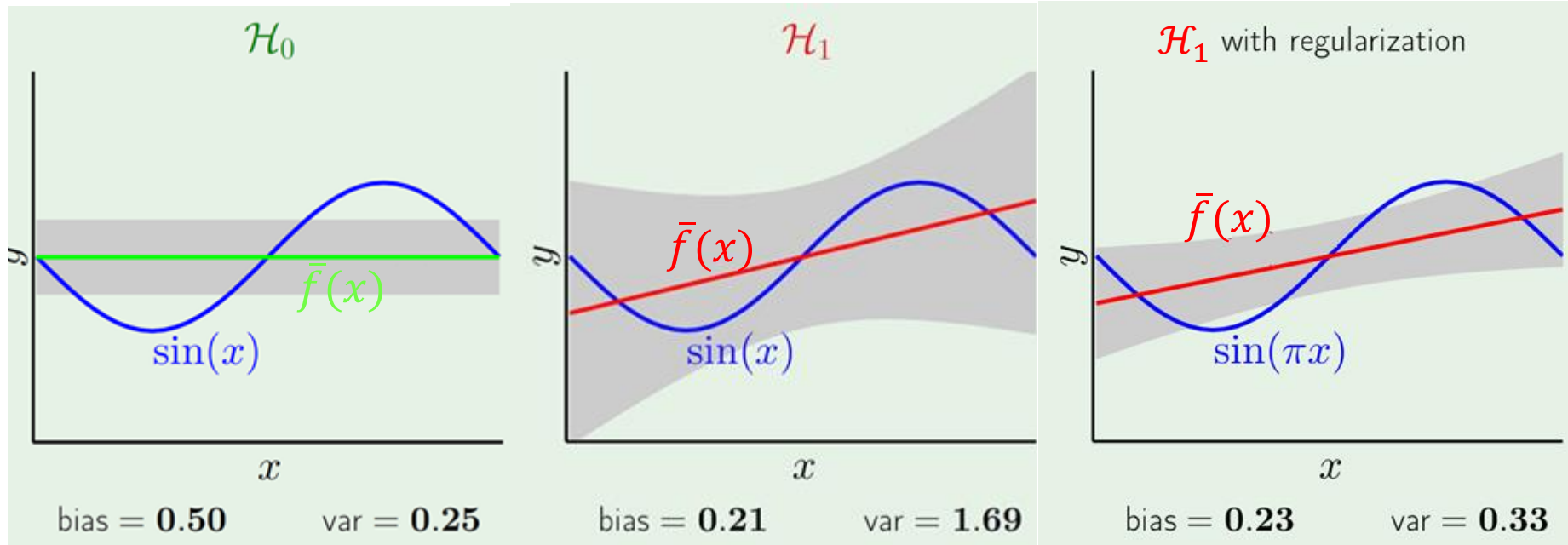


with regularization

Regularization: bias and variance

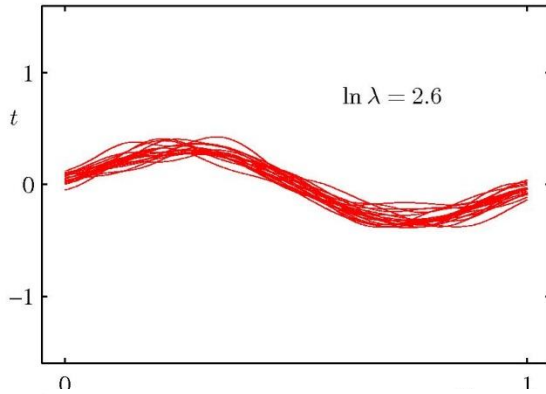


Winner of \mathcal{H}_0 , \mathcal{H}_1 , and \mathcal{H}_1 with regularization

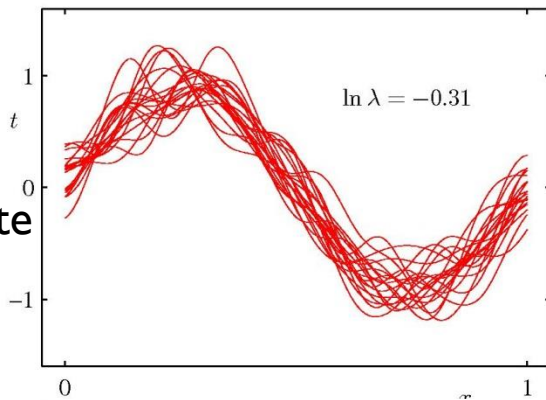


Regularization and bias/variance

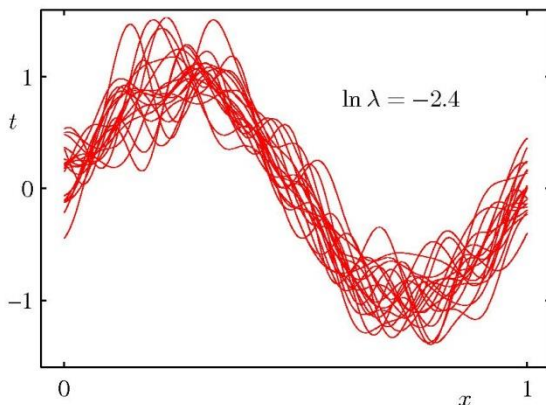
λ is large



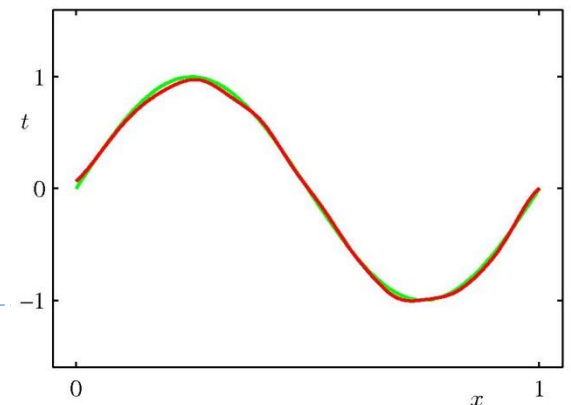
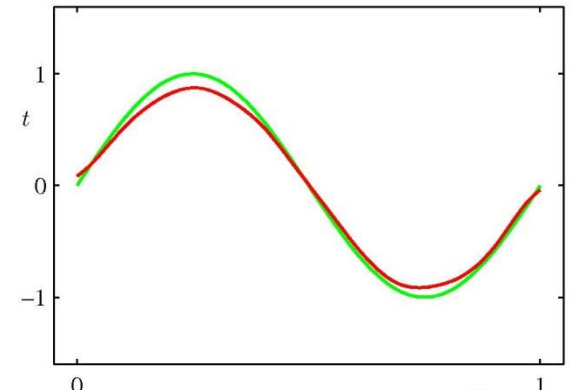
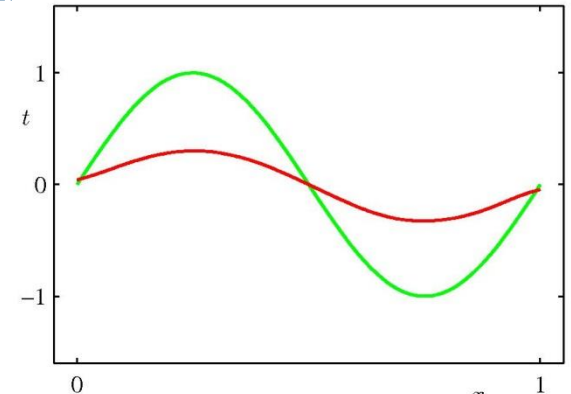
λ is intermediate



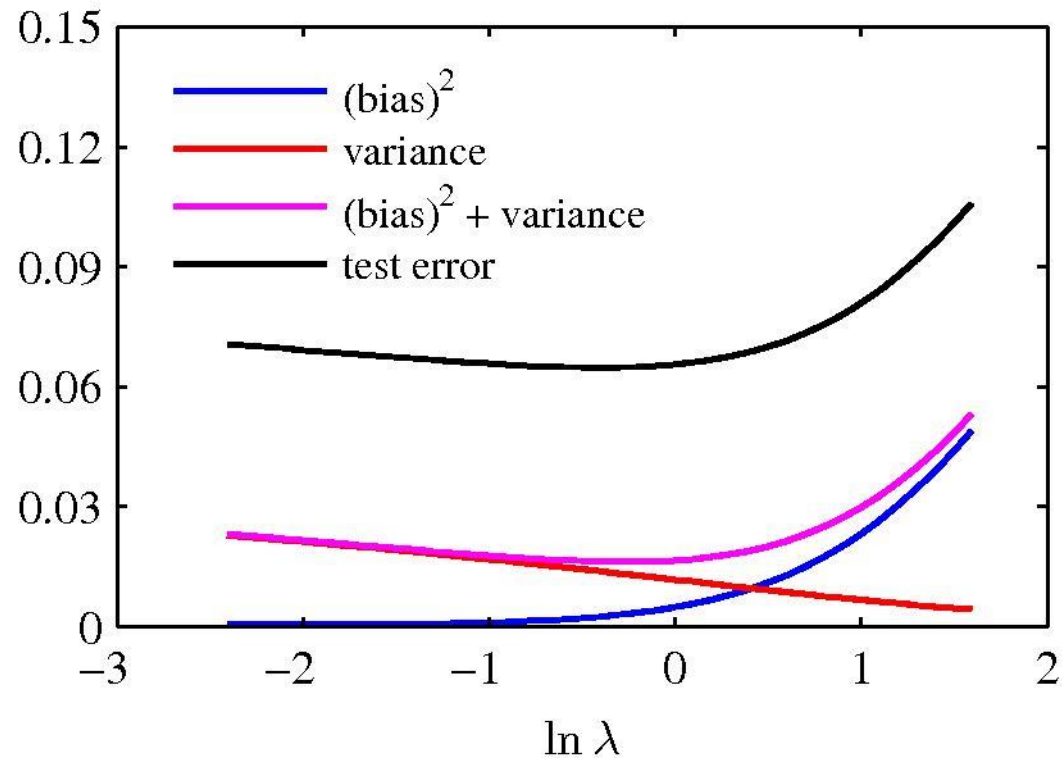
λ is small



$L = 100$ data sets
 $n = 25$
 $m = 25$



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 1.1, 1.3, 3.1, 3.2.
- ▶ Yaser S. Abu-Mostafa, Malik Maghdouh-Ismael, and Hsuan-Tien Lin, “**Learning from Data**”, Chapter 2.3, 3.2, 3.4.