

CSE 417T

Introduction to Machine Learning

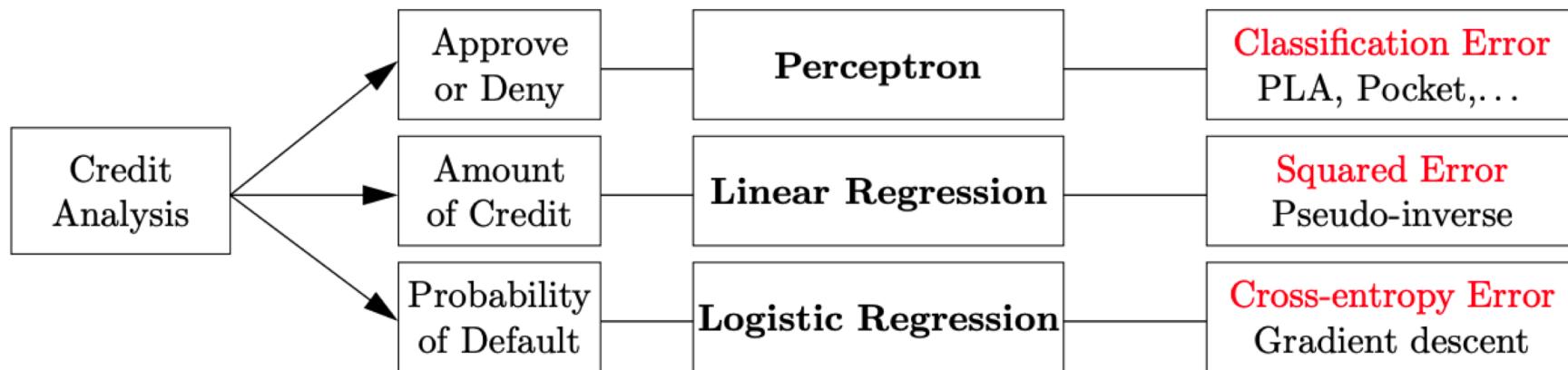
Lecture 10
Instructor: Chien-Ju (CJ) Ho

Recap

Linear Models

This is why it's called linear models

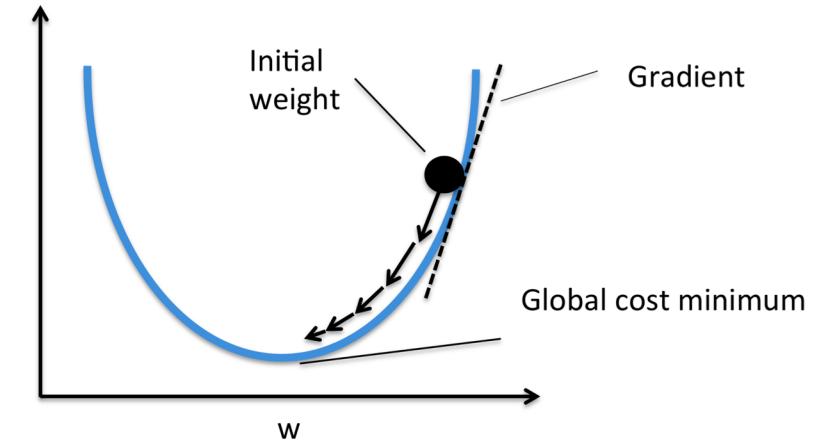
- H contains hypothesis $h(\vec{x})$ as **some function of $\vec{w}^T \vec{x}$**



- Algorithm:
 - Focus on $g = \operatorname{argmin}_{h \in H} E_{in}(h)$
 - Gradient descent is one of the common optimization algorithms

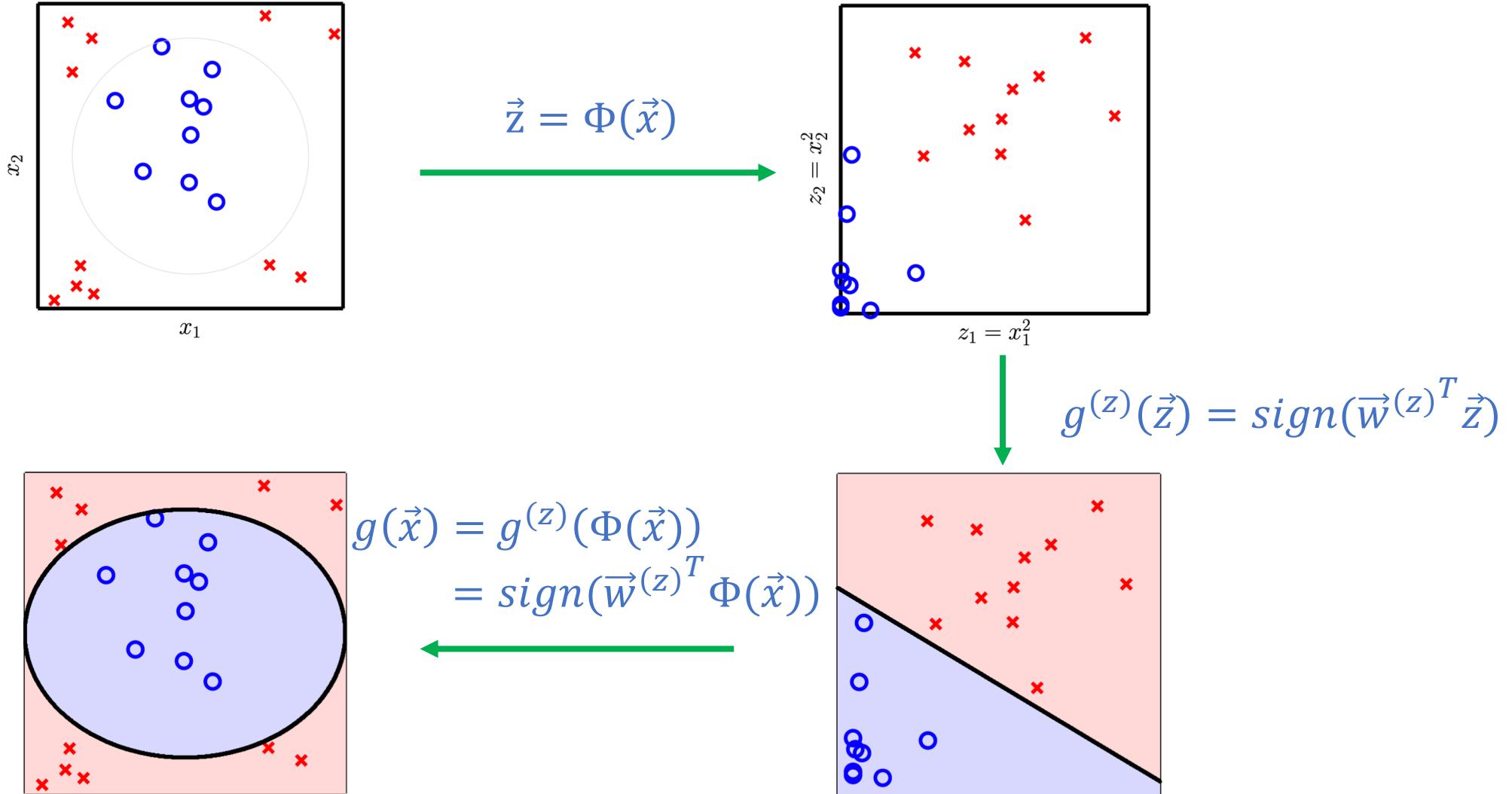
Gradient Descent

- Gradient descent algorithm
 - Initialize $\vec{w}(0)$
 - For $t = 0, \dots$
 - $\vec{w}(t + 1) \leftarrow \vec{w}(t) - \eta \nabla_{\vec{w}} E_{in}(\vec{w}(t))$
 - Terminate if the stop conditions are met
 - Return the final weights
- Stochastic gradient decent
 - Replace the update step:
 - Randomly choose n from $\{1, \dots, N\}$
 - $\vec{w}(t + 1) \leftarrow \vec{w}(t) - \eta \nabla_{\vec{w}} e_n(\vec{w}(t))$



Works for functions where gradient exists everywhere

Nonlinear Transformation



MUST Choose Φ BEFORE Looking at the Data

- Rely on domain knowledge (feature engineering)
 - Handwriting digit recognition example
- Use common sets of feature transformation
 - Polynomial transformation
 - E.g., 2nd order Polynomial transformation
 - $\vec{x} = (1, x_1, x_2)$, $\Phi_2(\vec{x}) = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$
 - Plus: more powerful (contains circle, ellipse, hyperbola, etc)
 - Minus:
 - More computation/storage
 - Worse generalization error

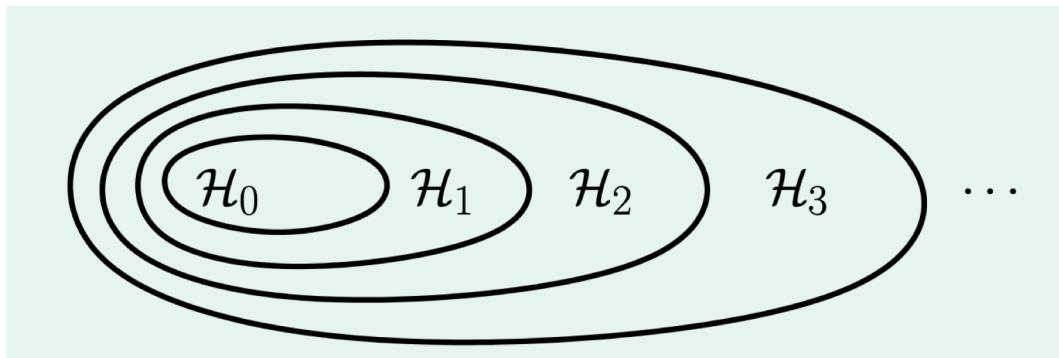
The VC dimension of d-dim perceptron is d+1

Q-th Order Polynomial Transform

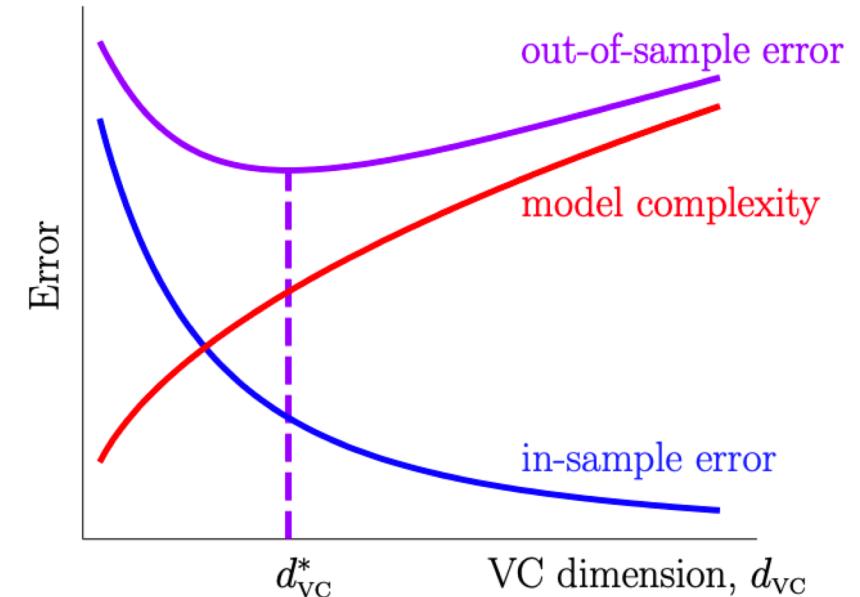
- $\vec{x} = (1, x_1, \dots, x_d)$
- $\Phi_1(\vec{x}) = \vec{x}$
- $\Phi_Q(\vec{x}) = (\Phi_{Q-1}(\vec{x}), x_1^Q, x_1^{Q-1}x_2, \dots, x_d^Q)$
- Each element in $\Phi_Q(\vec{x})$ is in the form of $\sum_{i=1}^d x_i^{a_i}$
 - where $\sum_{i=1}^d a_i \leq Q$, and a_i is a non-negative integer
- Number of elements in $\Phi_Q(\vec{x})$: $\binom{Q+d}{Q}$ (including the initial 1)

Structural Hypothesis Sets

- Let H_Q be the linear model for the $\Phi_Q(\vec{x})$ space



- Let $g_Q = \operatorname{argmin}_{h \in H_Q} E_{in}(h)$
 - $H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots$
 - $d_{vc}(H_1) \leq d_{vc}(H_2) \leq \dots$
 - $E_{in}(g_1) \geq E_{in}(g_2) \geq \dots$



Brief Lecture Notes Today

The notes are not intended to be comprehensive. They should be accompanied by lectures and/or textbook.
Let me know if you spot errors.

Overfitting

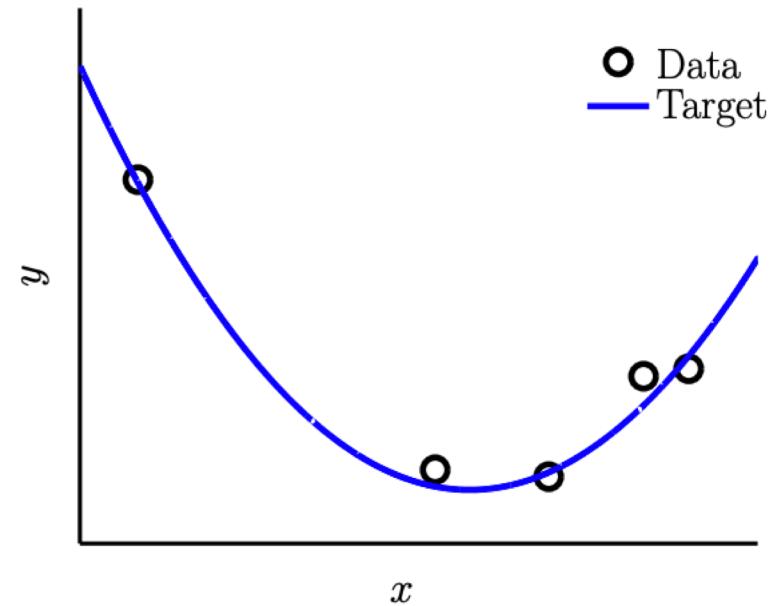
[Adapted from the slides by Malik Magdon-Ismail]

Setup of the Discussion

- Regression with polynomial transform
 - Input: 1-dimensional x
 - $\Phi_Q(x) = (1, x, x^2, x^3, \dots, x^Q)$
 - $H_Q = \{h(x) = w_0 + w_1x + w_2x^2 + \dots + w_Qx^Q\}$
- Q th-order polynomial fit
 - Solve linear regression on the $\Phi_Q(\vec{x})$ space using H_Q
 - Looking to minimize E_{in} : $g_Q = \operatorname{argmin}_{h \in H_Q} E_{in}(h)$

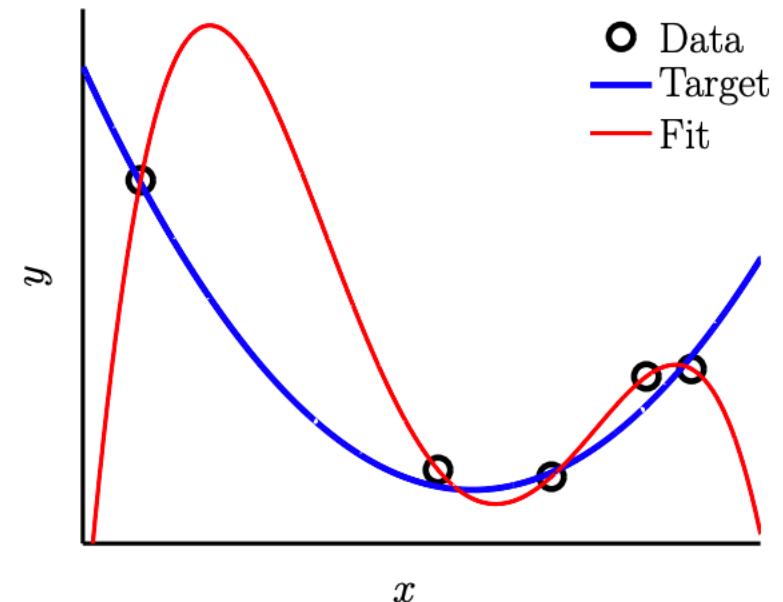
A Simple Example

- Target f : 4th order function
- # data points: $N = 5$
- Small noise:
 - $y = f(x) + \epsilon$ with small ϵ
- 4th order polynomial fit
 - $h(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$
 - Find $g_4 = \text{argmin}_h E_{in}(h)$



A Simple Example

- Target f : 4th order function
- # data points: $N = 5$
- Small noise:
 - $y = f(x) + \epsilon$ with small ϵ
- 4th order polynomial fit
 - $h(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$
 - Find $g_4 = \text{argmin}_h E_{in}(h)$



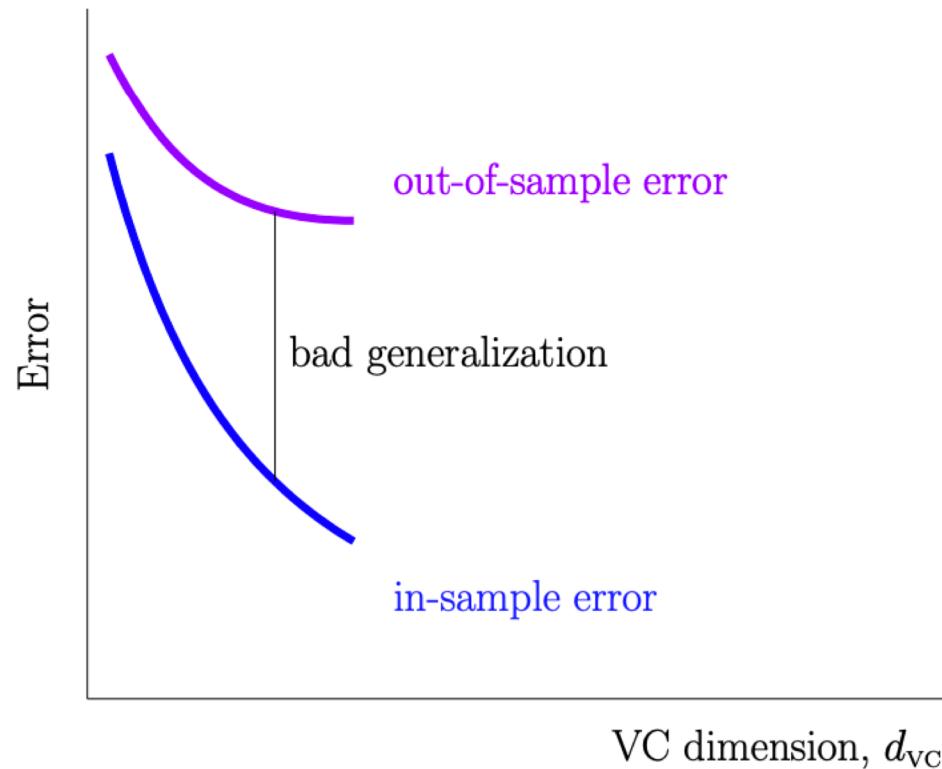
Classical overfitting: $E_{in} = 0$, but lead to a large E_{out}

Fitting the **noise** instead of the data

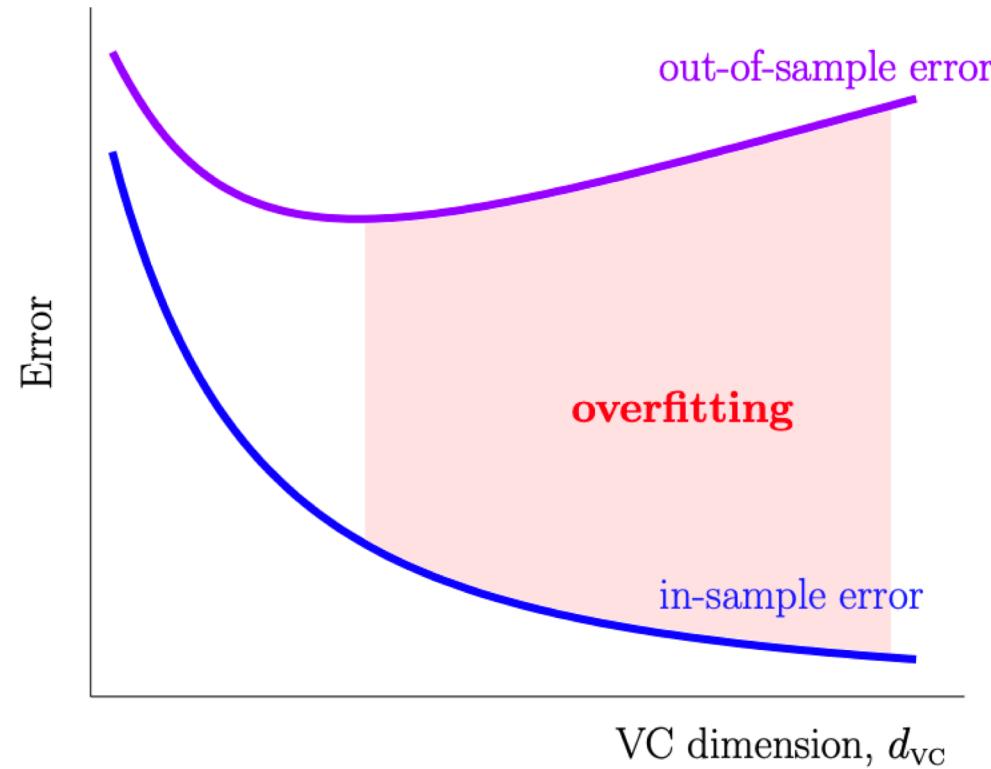
What is Overfitting?

Fitting the data **more** than is **warranted**

Overfitting is Not Just Bad Generalization



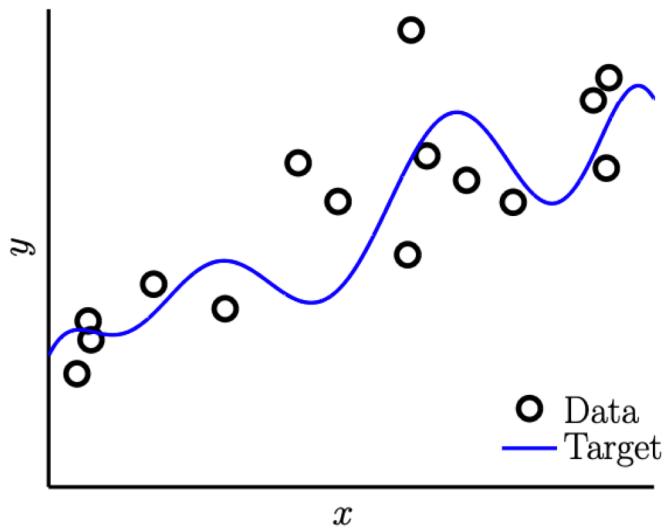
Overfitting is Not Just Bad Generalization



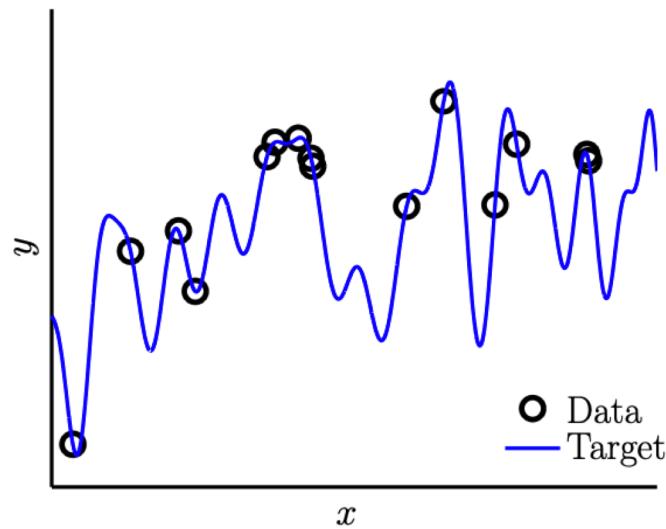
Overfitting

Going for lower and lower E_{in} results in higher and higher E_{out}

Case Study: 2nd vs 10th Order Polynomial Fit



10th order f with noise.



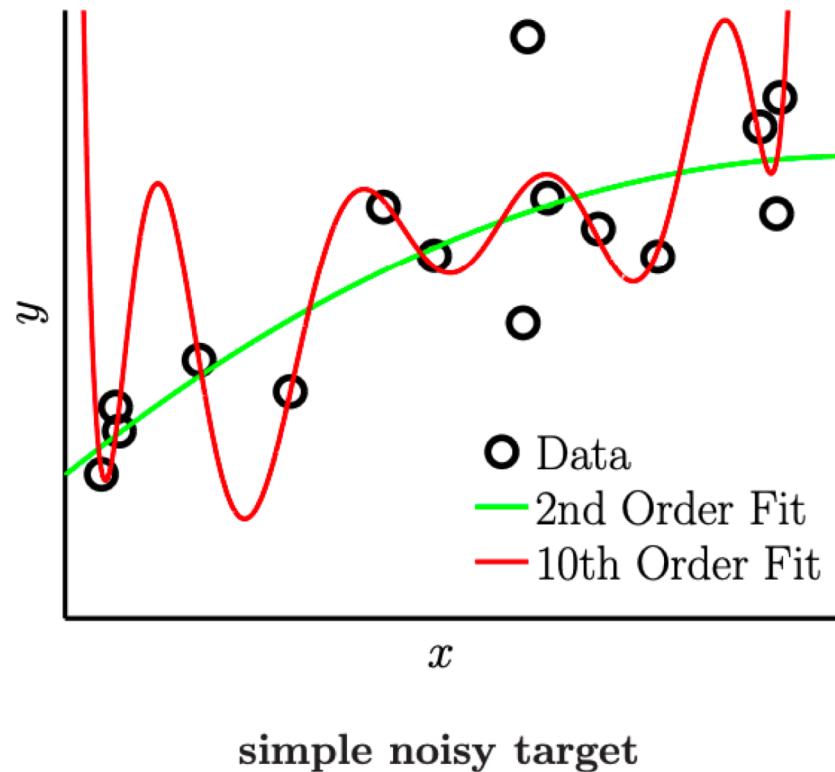
50th order f with no noise.

H_2 : 2nd order polynomial fit

H_{10} : 10th order polynomial fit

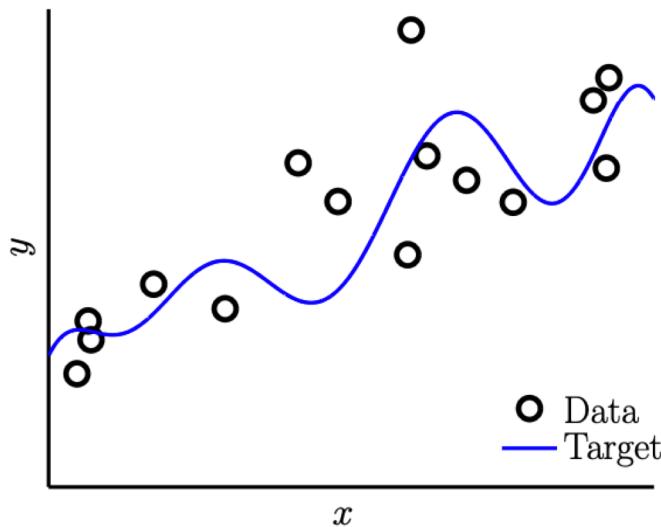
Which model do you choose for the left problem and why?

Target Function: 10th Order f with Noise

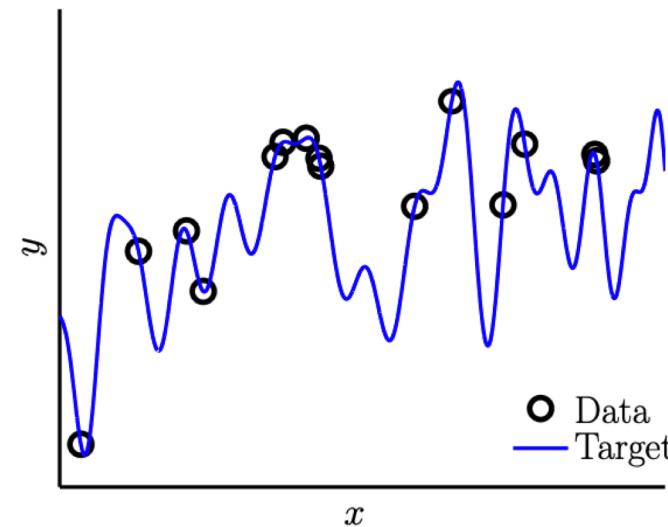


	2nd Order	10th Order
E_{in}	0.050	0.034
E_{out}	0.127	9.00

- Irony of two learners O and R
- Both know the target is 10th order
- O chooses H_{10}
- R chooses H_2
- R outperforms O



10th order f with noise.



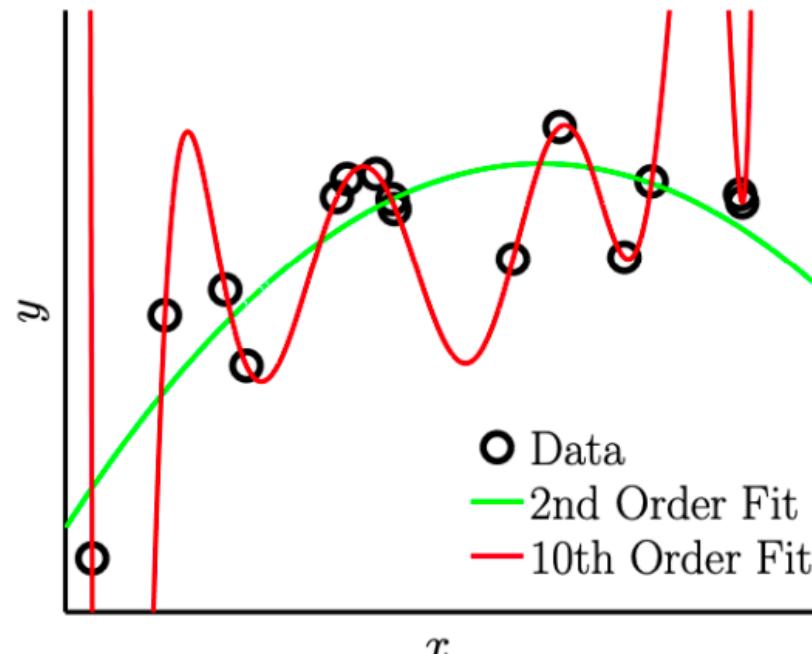
50th order f with no noise.

H_2 : 2nd order polynomial fit

H_{10} : 10th order polynomial fit

Which model do you choose for the right problem and why?

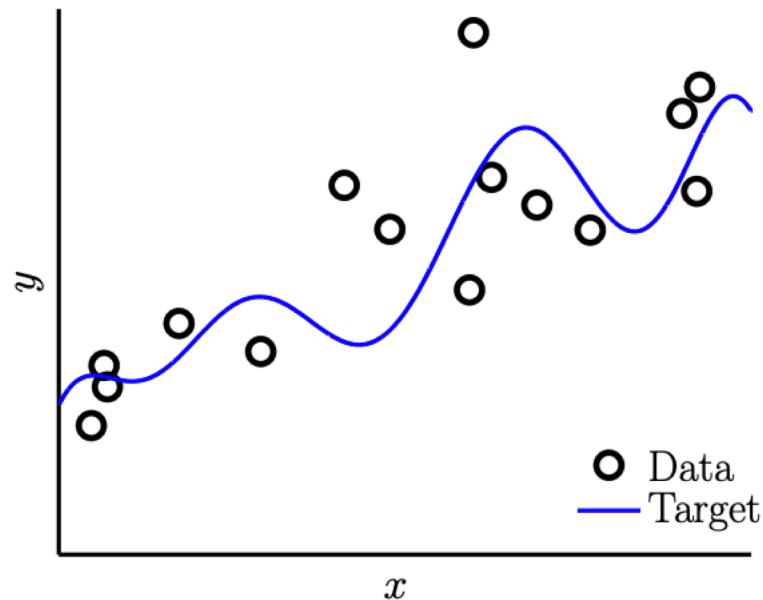
Simpler H is better even for complex target with **no noise**



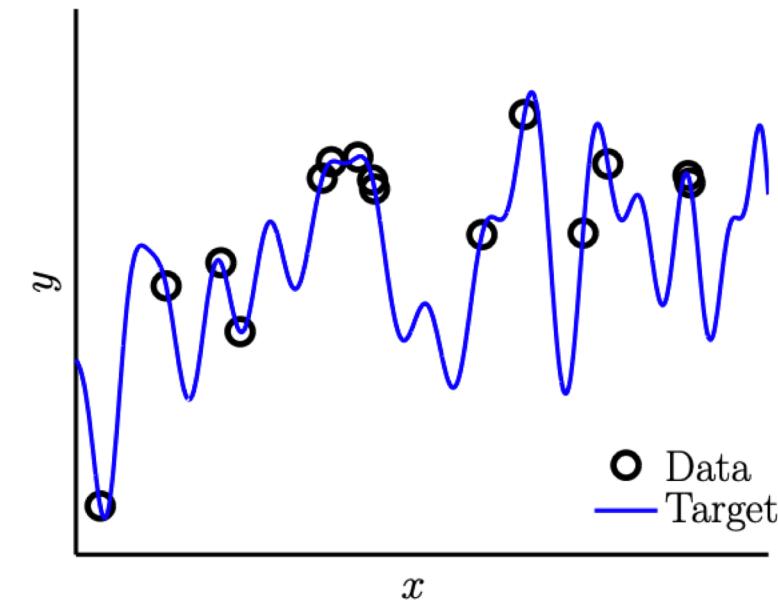
complex noiseless target

	2nd Order	10th Order
E_{in}	0.029	10^{-5}
E_{out}	0.120	7680

Is There Really “No Noise”?

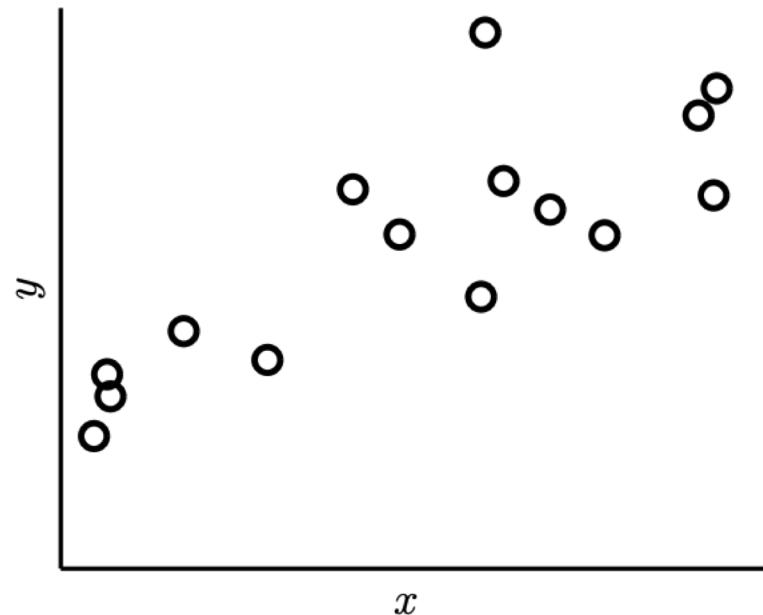


Simple f with noise.

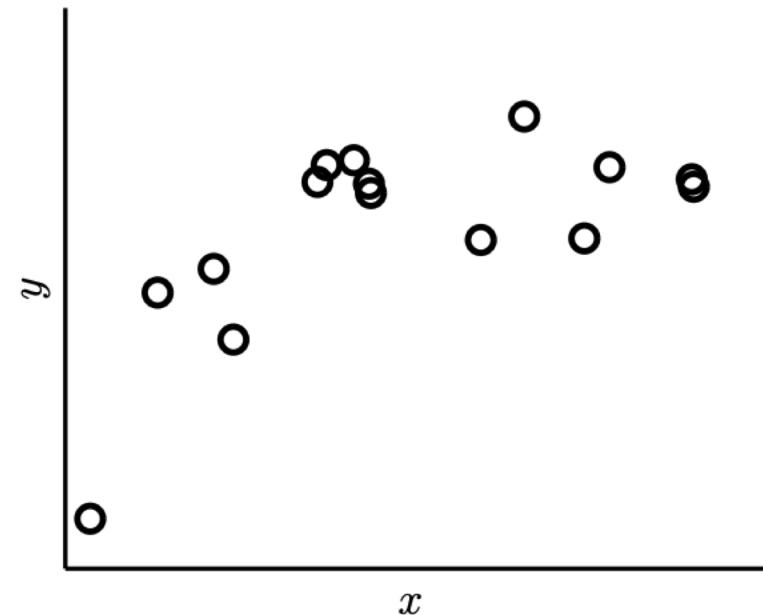


Complex f with no noise.

Is There Really “No Noise”?



Simple f with noise.

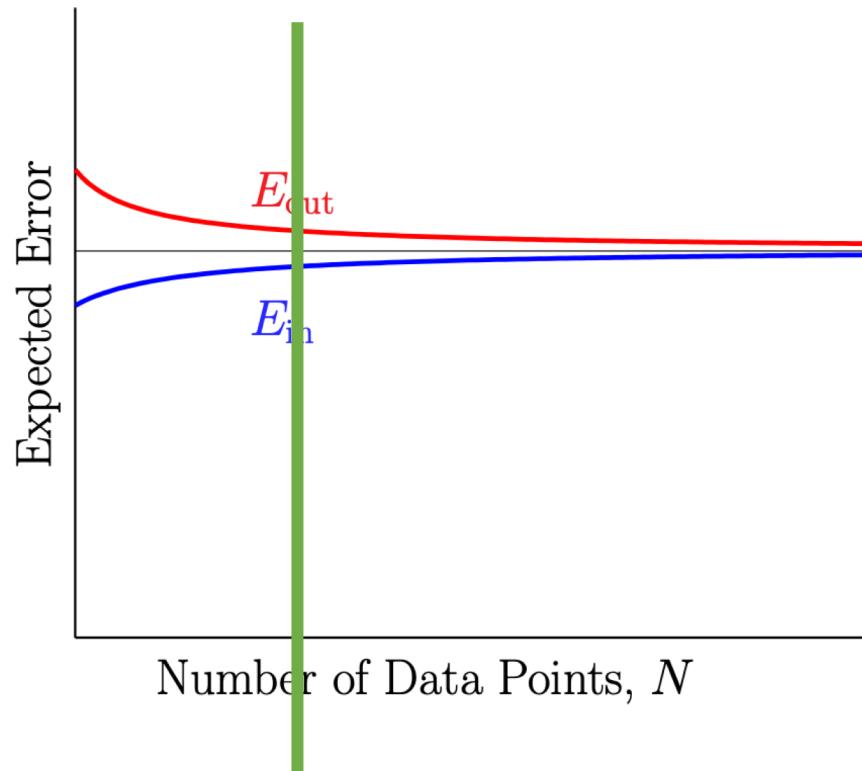


Complex f with no noise.

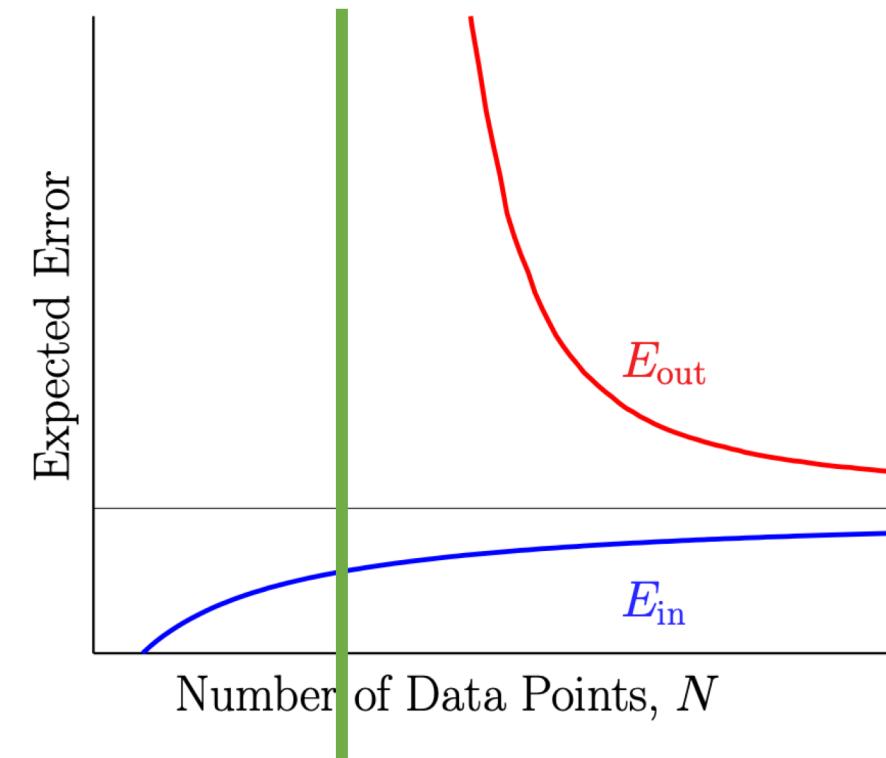
H should match quantity and quality of data, not (unknown) f

Why is H_2 Better than H_{10} ?

Learning curve for H_2



Learning curve for H_{10}



When N is small, $E_{out}(g_{10}) \geq E_{out}(g_2)$

A Detailed Experiment

Study the **level of noise** and **target complexity**, and **# data points N**

$$y = f(x) + \epsilon(x) = \sum_{q=0}^{Q_f} \alpha_q x^q + \epsilon(x)$$

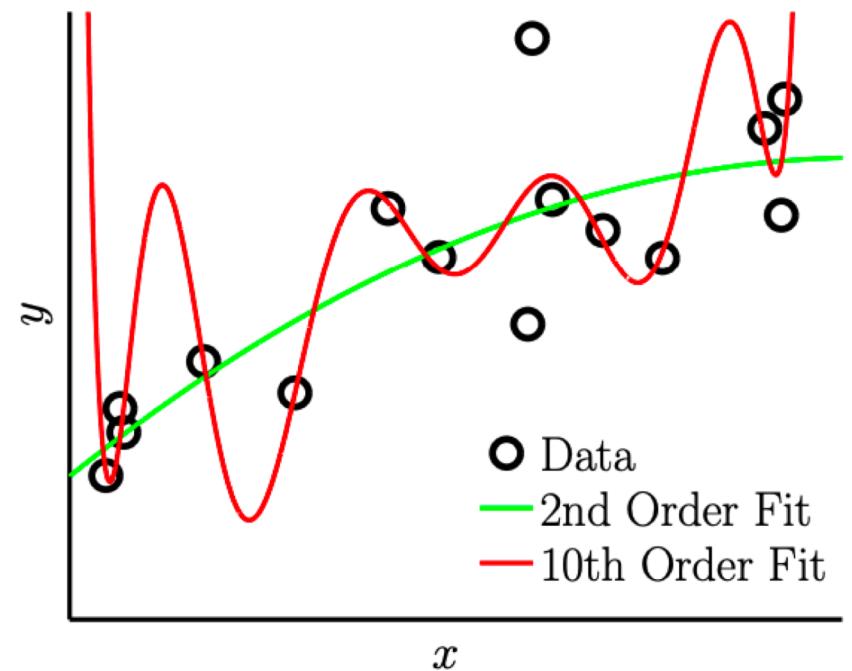
Noise level: variance σ^2 of $\epsilon(x)$

Target complexity: Q_f

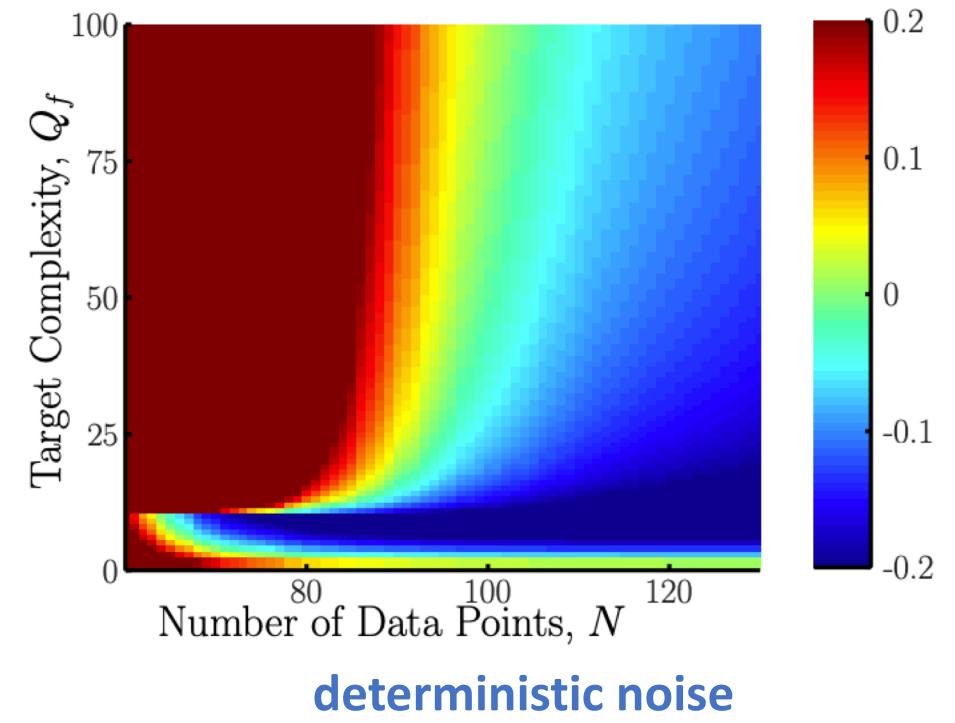
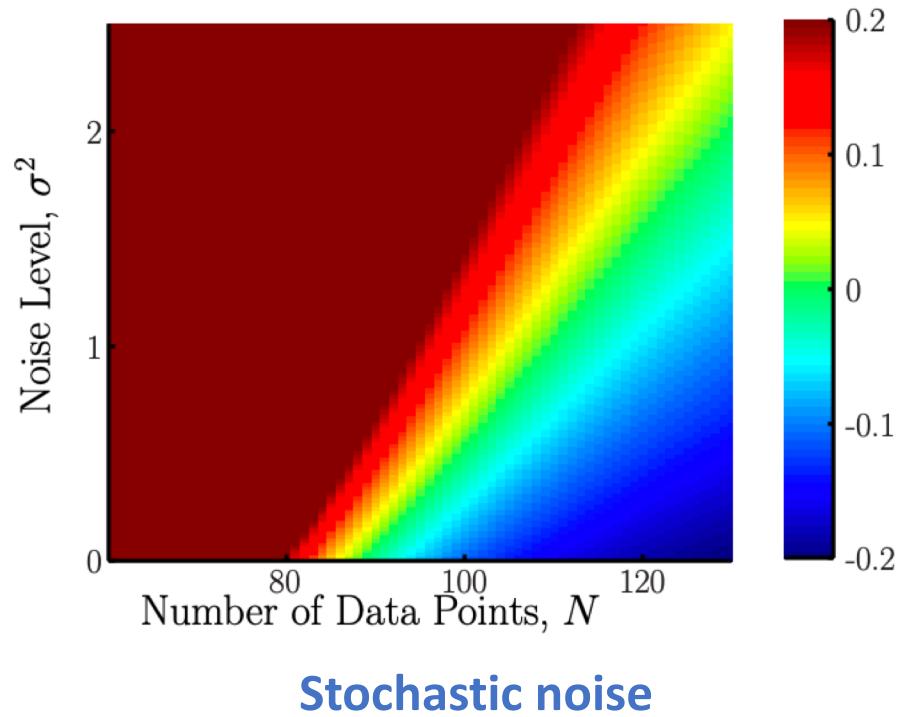
Data set size: N

The Overfit Measure

- Fit the data set using H_2 and H_{10}
 - Let g_2 and g_{10} be the learned hypothesis
- Overfit measure
 - $E_{out}(g_{10}) - E_{out}(g_2)$
 - This value is large if overfit happens



Overfit Measure: $E_{out}(g_{10}) - E_{out}(g_2)$



Number of data points \uparrow	Overfitting \downarrow
Noise \uparrow	Overfitting \uparrow
Target complexity \uparrow	Overfitting \uparrow

Noise:
The part of y we cannot model

Stochastic Noise

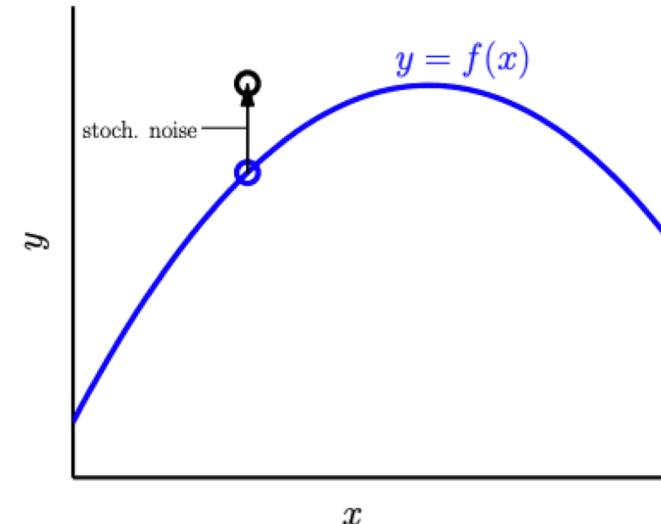
We would like to learn from \bullet :

$$y_n = f(x_n)$$

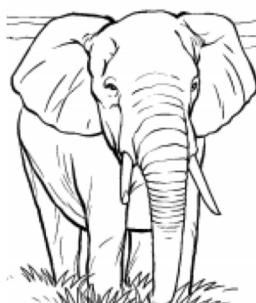
Unfortunately, we only observe \circ :

$$y_n = f(x_n) + \text{'stochastic noise'}$$

↑
no one can model this



Stochastic Noise: fluctuations/measurement errors we cannot model.



Deterministic Noise

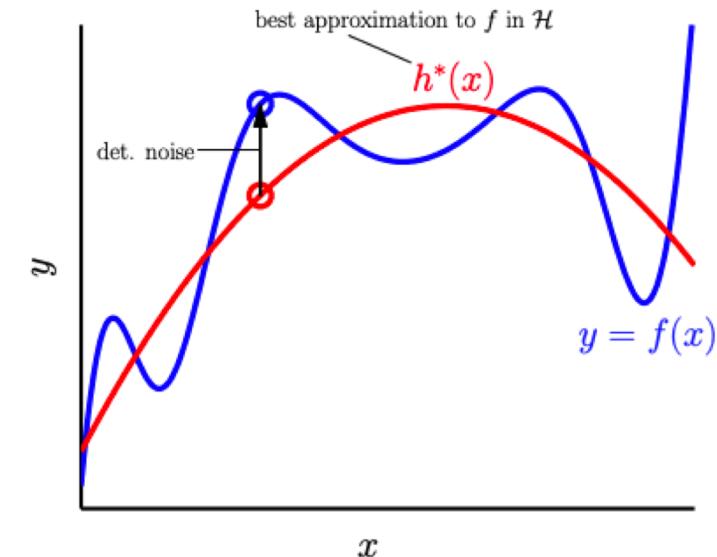
We would like to learn from $\textcolor{red}{O}$:

$$y_n = \textcolor{red}{h}^*(x_n)$$

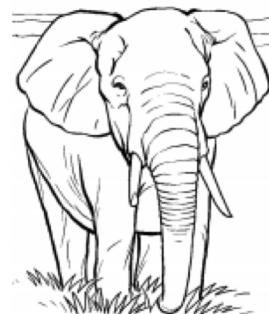
Unfortunately, we only observe $\textcolor{blue}{O}$:

$$\begin{aligned} y_n &= \textcolor{blue}{f}(x_n) \\ &= \textcolor{red}{h}^*(x_n) + \text{'deterministic noise'} \end{aligned}$$

\uparrow
 \mathcal{H} cannot model this



Deterministic Noise: the part of f we cannot model.



Deterministic Noise

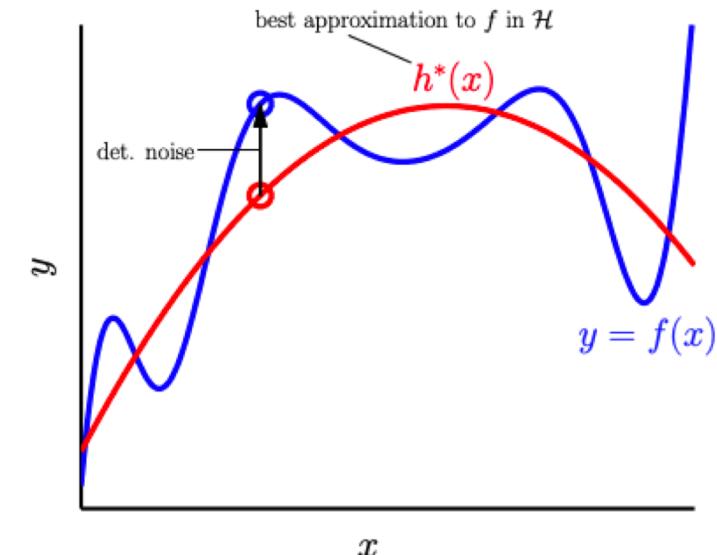
We would like to learn from $\textcolor{red}{O}$:

$$y_n = \textcolor{red}{h}^*(x_n)$$

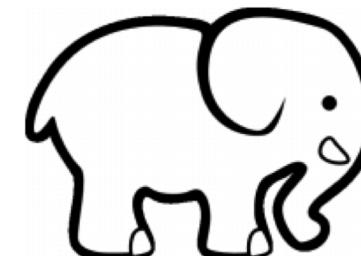
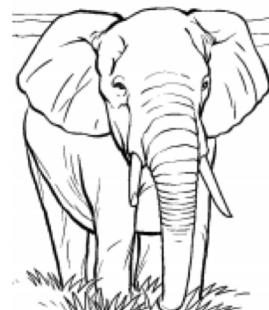
Unfortunately, we only observe $\textcolor{blue}{O}$:

$$\begin{aligned} y_n &= \textcolor{blue}{f}(x_n) \\ &= \textcolor{red}{h}^*(x_n) + \text{'deterministic noise'} \end{aligned}$$

\uparrow
 \mathcal{H} cannot model this

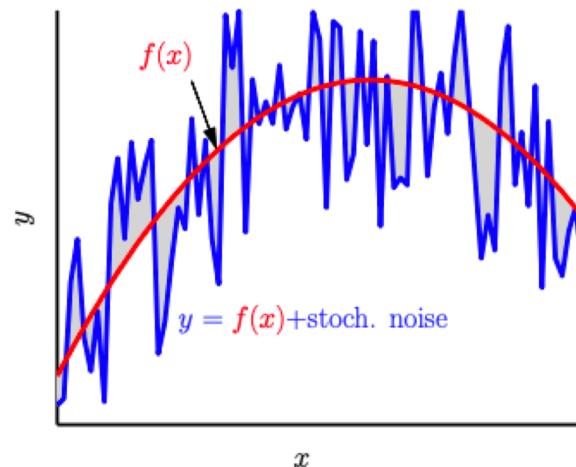


Deterministic Noise: the part of f we cannot model.



Both sources of noises hurt learning

Stochastic Noise



source: random measurement errors

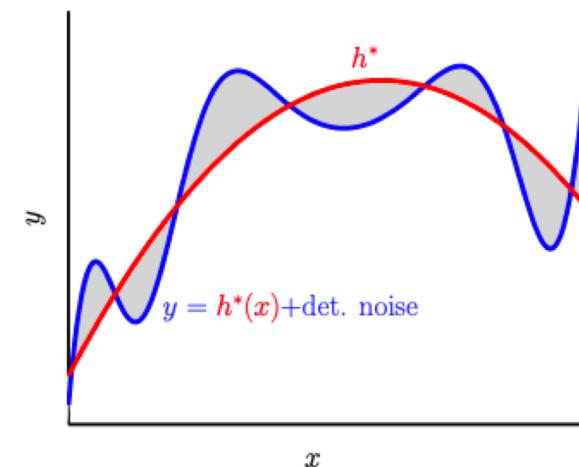
re-measure y_n

stochastic noise changes.

change \mathcal{H}

stochastic noise the same.

Deterministic Noise



source: learner's \mathcal{H} cannot model f

re-measure y_n

deterministic noise the same.

change \mathcal{H}

deterministic noise changes.

We have single \mathcal{D} and fixed \mathcal{H} so we cannot distinguish

Noise and Bias-Variance Decomposition

$$y = f(\vec{x}) + \epsilon$$

$$\mathbb{E}[E_{out}(\vec{x})] = \sigma^2 + \text{bias} + \text{variance}$$



Stochastic Noise

Deterministic noise

How to Fight Overfitting

- Regularization
- Validation
- (The focus of the next two lectures)