

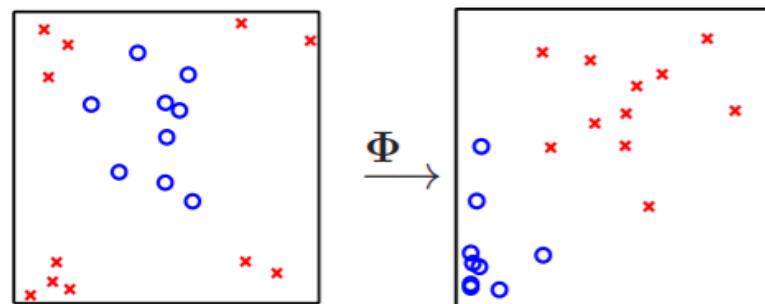
# Machine Learning from Data

Lecture 11: Spring 2021

# Today's Lecture

- Overfitting
  - What is overfitting?
  - When does it occur?
  - Stochastic Vs. Deterministic Noise

# Non-Linear Transforms

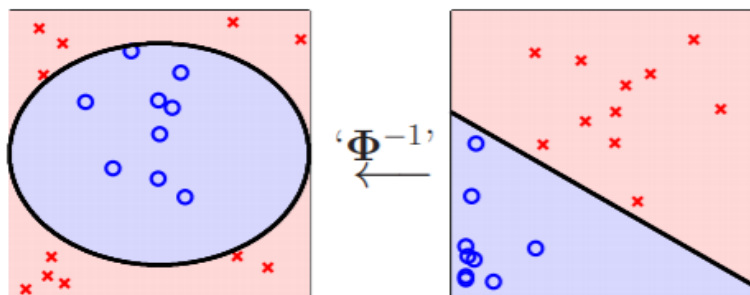


1. Original data

$$\mathbf{x}_n \in \mathcal{X}$$

2. Transform the data

$$\mathbf{z}_n = \Phi(\mathbf{x}_n) \in \mathcal{Z}$$



4. Classify in  $\mathcal{X}$ -space

$$g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}))$$

3. Separate data in  $\mathcal{Z}$ -space

$$\tilde{g}(\mathbf{z}) = \text{sign}(\tilde{\mathbf{w}}^T \mathbf{z})$$

$\mathcal{X}$ -space is  $\mathbb{R}^d$

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$$

$$y_1, y_2, \dots, y_N$$

no weights

$$d_{\text{vc}} = d + 1$$

$$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}))$$

$\mathcal{Z}$ -space is  $\mathbb{R}^{\bar{d}}$

$$\mathbf{z} = \Phi(\mathbf{x}) = \begin{bmatrix} 1 \\ \Phi_1(\mathbf{x}) \\ \vdots \\ \Phi_{\bar{d}}(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 1 \\ z_1 \\ \vdots \\ z_{\bar{d}} \end{bmatrix}$$

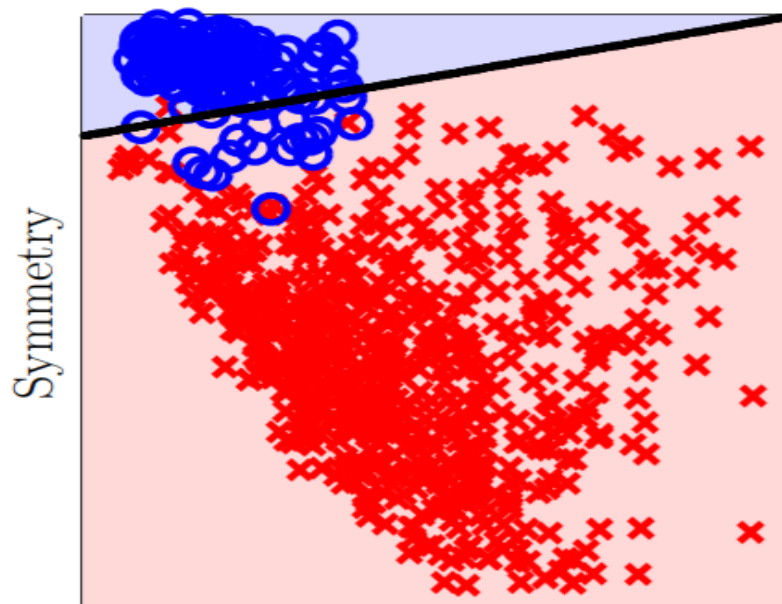
$$\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$$

$$y_1, y_2, \dots, y_N$$

$$\tilde{\mathbf{w}} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{\bar{d}} \end{bmatrix}$$

$$d_{\text{vc}} = \bar{d} + 1$$

# Digits Data

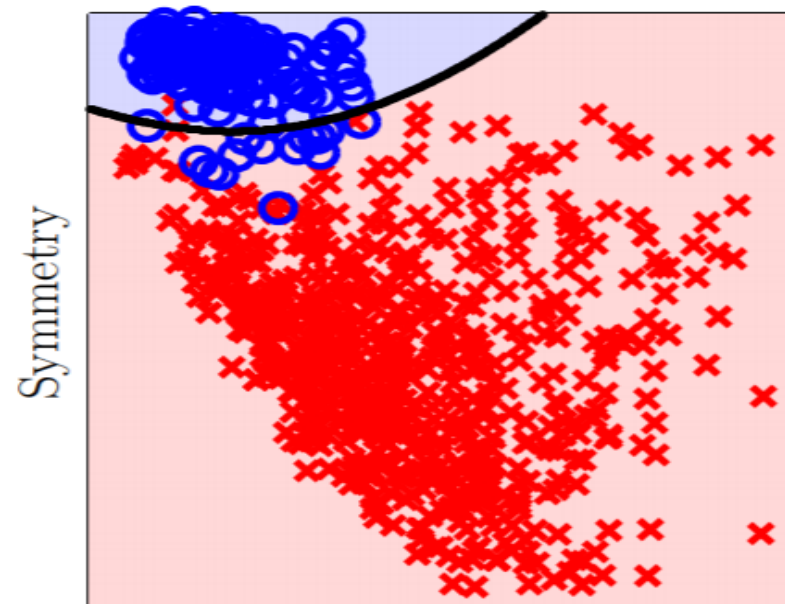


Average Intensity

**Linear model**

$$E_{\text{in}} = 2.13\%$$

$$E_{\text{out}} = 2.38\%$$



Average Intensity

**3rd order polynomial model**

$$E_{\text{in}} = 1.75\%$$

$$E_{\text{out}} = 1.87\%$$

# Humans Overfit (Superstitions)

- Fear of Friday the 13<sup>th</sup>

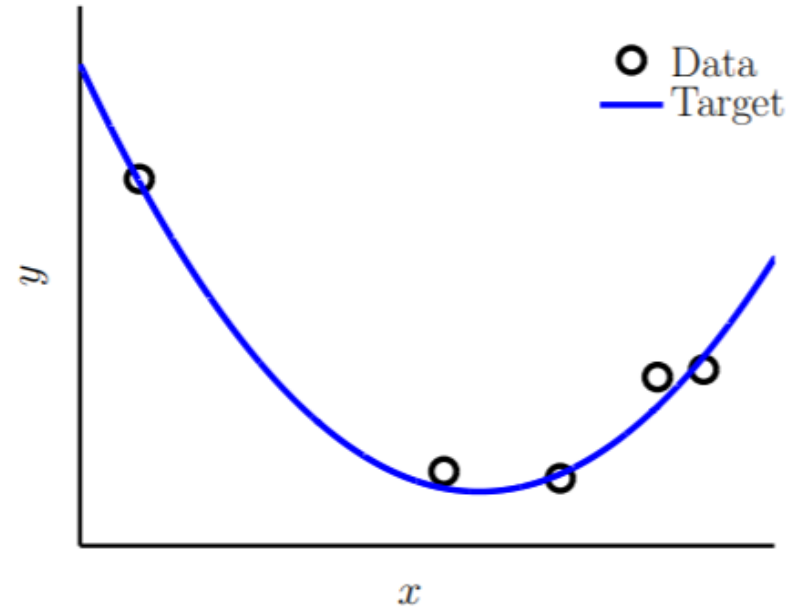
# Illustration of Overfitting

Quadratic  $f$

5 data points

A *little* noise (measurement error)

5 data points  $\rightarrow$  4th order polynomial fit



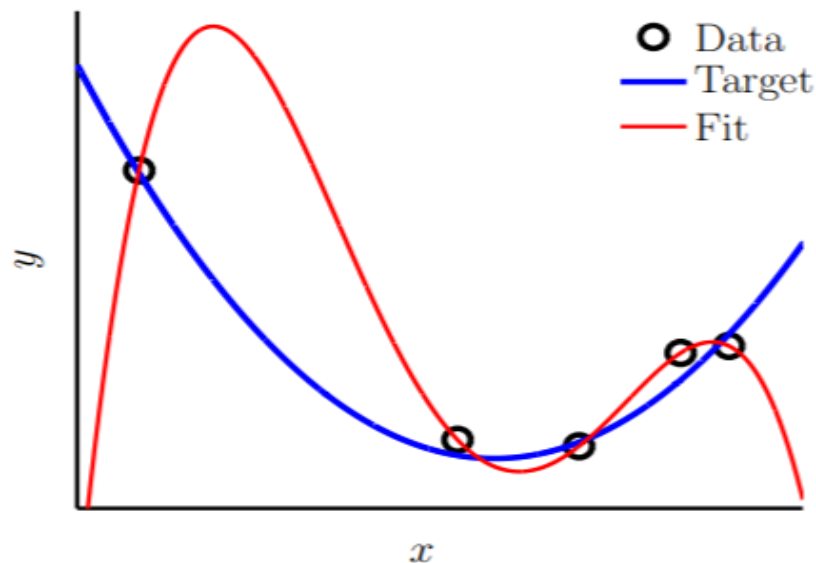
# Overfitting Example

Quadratic  $f$

5 data points

A *little* noise (measurement error)

5 data points  $\rightarrow$  4th order polynomial fit



Classic overfitting: simple target with excessively complex  $\mathcal{H}$ .

$$E_{\text{in}} \approx 0; E_{\text{out}} \gg 0$$

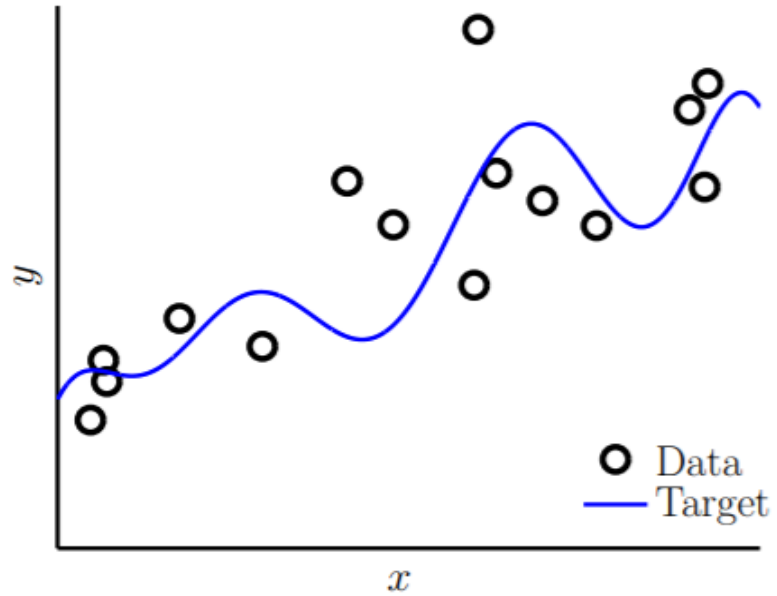
# Define Overfitting

- Fitting the data more than is warranted.

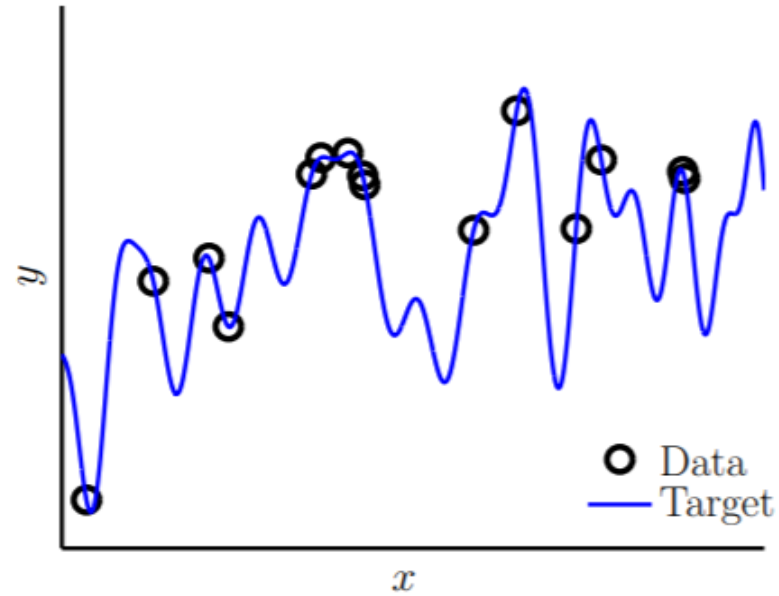




# Case Study



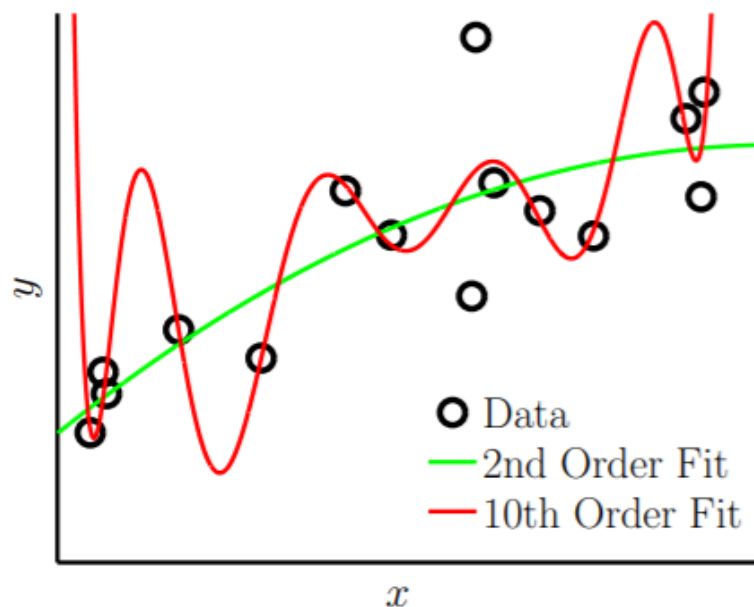
10th order  $f$  with noise.



50th order  $f$  with no noise.

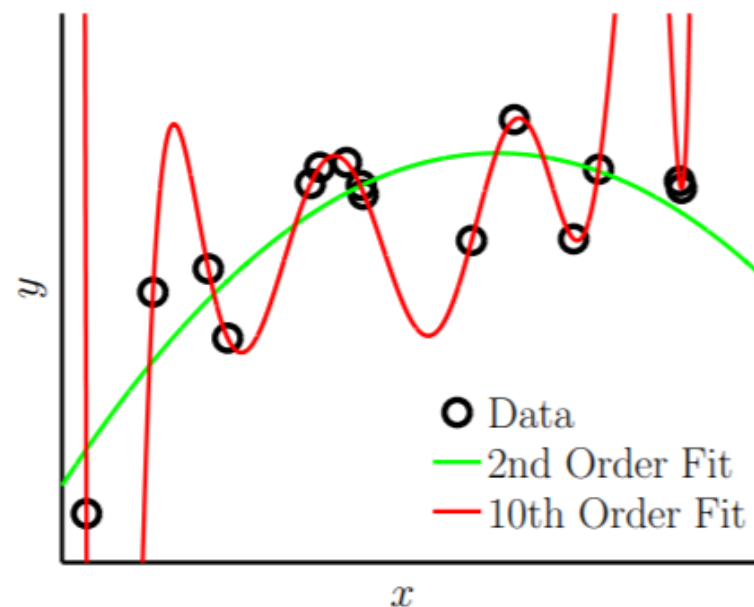


# 2<sup>nd</sup> order vs. 10<sup>th</sup> order polynomial



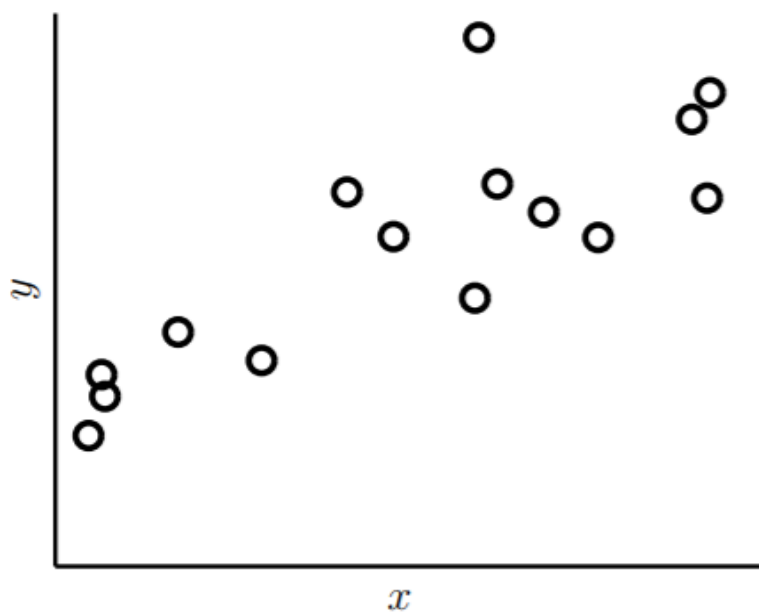
simple noisy target

	2nd Order	10th Order
$E_{\text{in}}$	0.050	0.034
$E_{\text{out}}$	0.127	<b>9.00</b>

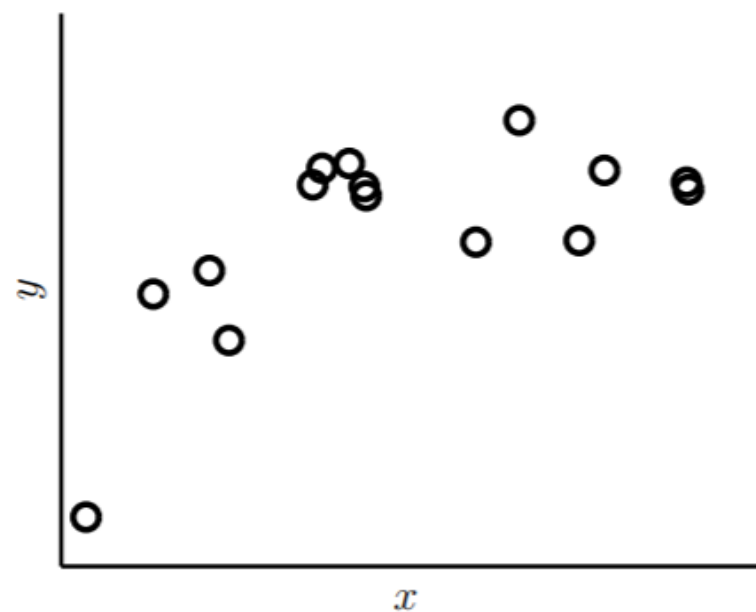


complex noiseless target

	2nd Order	10th Order
$E_{\text{in}}$	0.029	$10^{-5}$
$E_{\text{out}}$	0.120	<b>7680</b>



Simple  $f$  with noise.



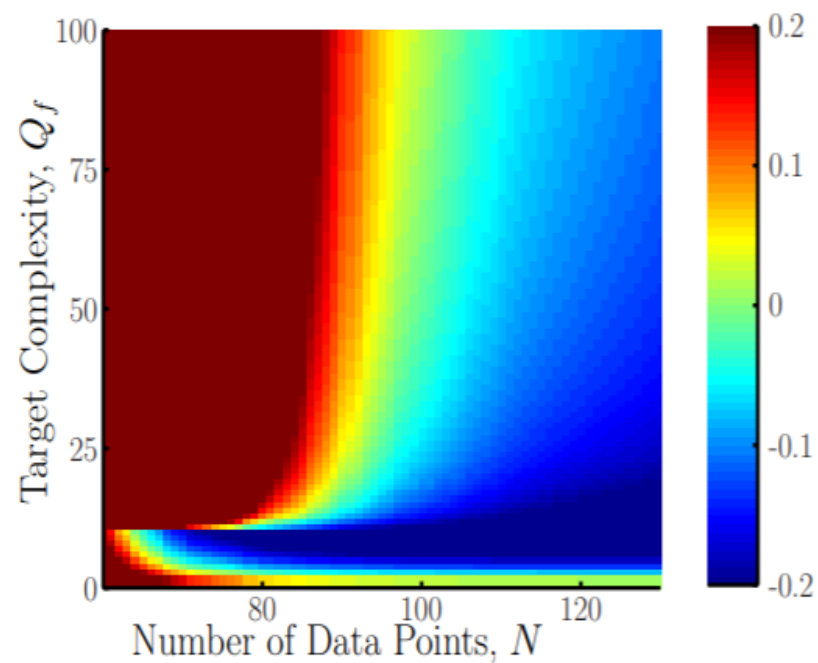
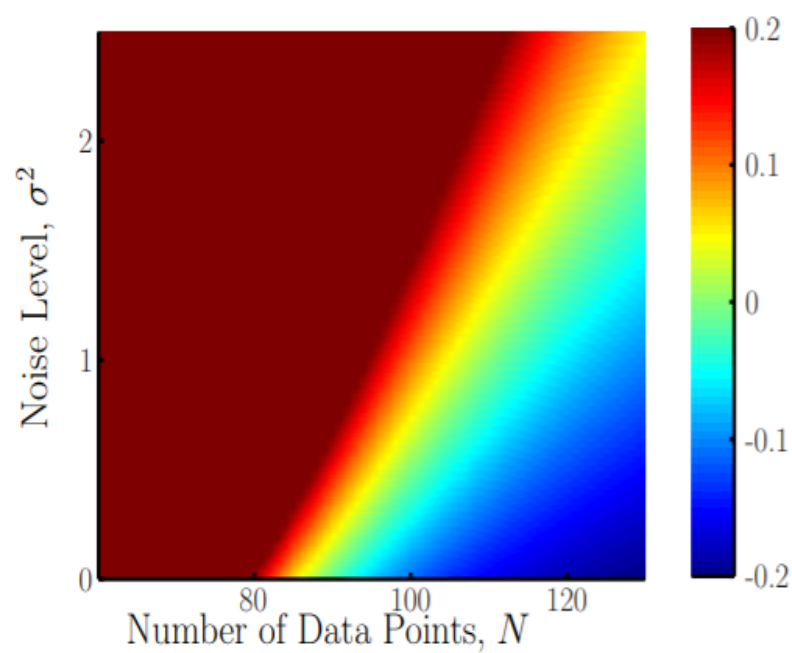
Complex  $f$  with no noise.

$\mathcal{H}$  should match *quantity and quality of data*, not  $f$

# Measure Overfitting

-

## Overfit Measure: $E_{\text{out}}(\mathcal{H}_{10}) - E_{\text{out}}(\mathcal{H}_2)$



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

# Noise

-





# Stochastic Noise

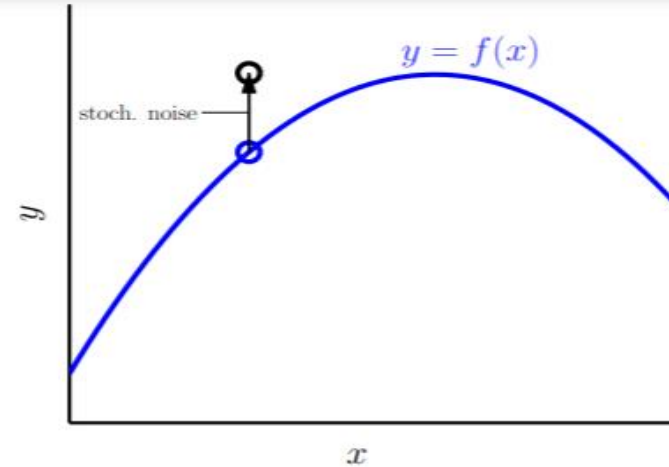
We would like to learn from ○:

$$y_n = f(x_n)$$

Unfortunately, we only observe ○:

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

↑  
no one can model this



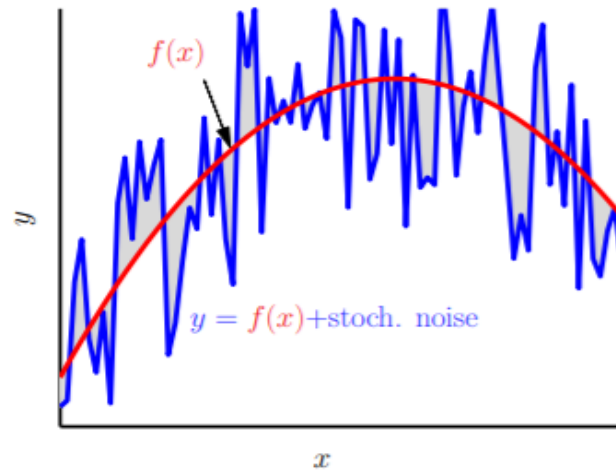
**Stochastic Noise:** fluctuations/measurement errors we cannot model.







### Stochastic Noise



**source:** random measurement errors

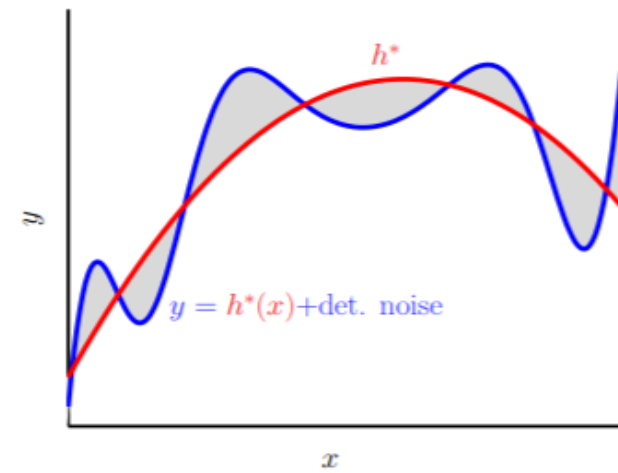
re-measure  $y_n$

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

# Deterministic Noise

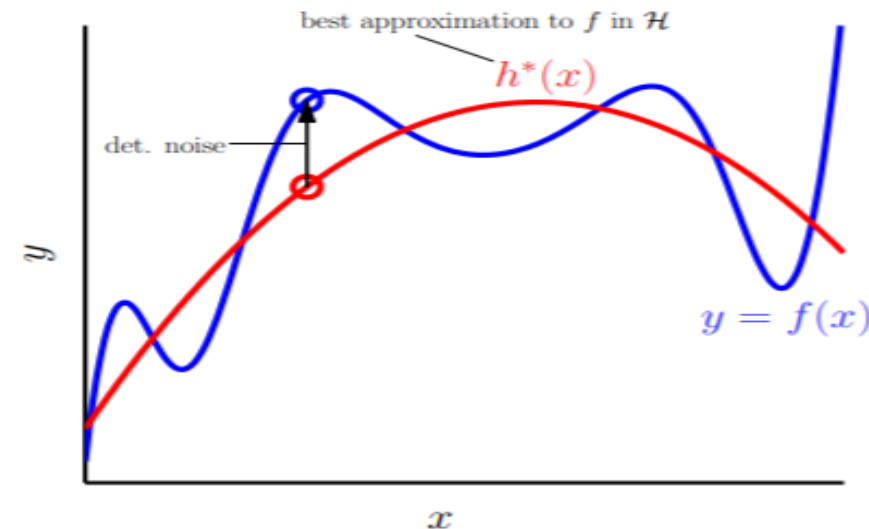
We would like to learn from  $\circ$ :

$$y_n = h^*(x_n)$$

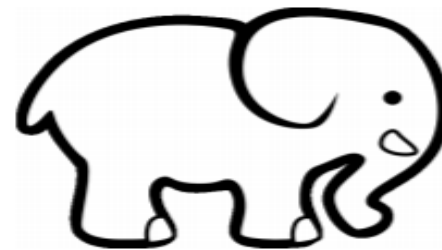
Unfortunately, we only observe  $\circ$ :

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

↑  
 $\mathcal{H}$  cannot model this



**Deterministic Noise:** the part of  $f$  we cannot model.



# Bias Variance Analysis and Noise

-







# Summary

-

Thanks!