

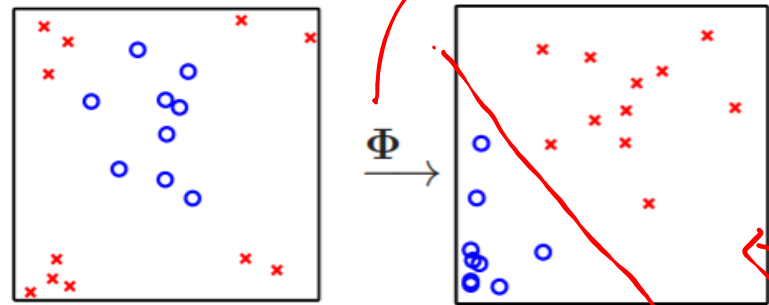
# Machine Learning from Data

Lecture 11: Spring 2021

# Today's Lecture

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

# Non-Linear Transforms

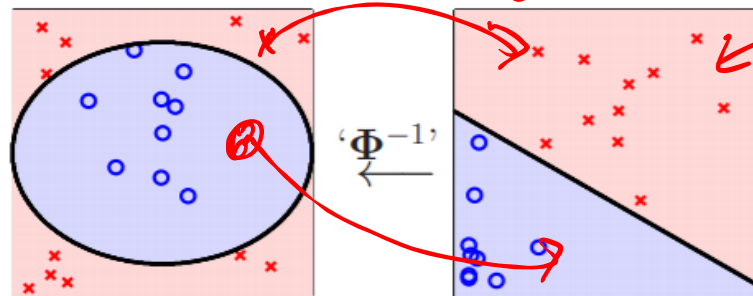


1. Original data

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

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  $\mathbb{R}^{\bar{d}}$

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

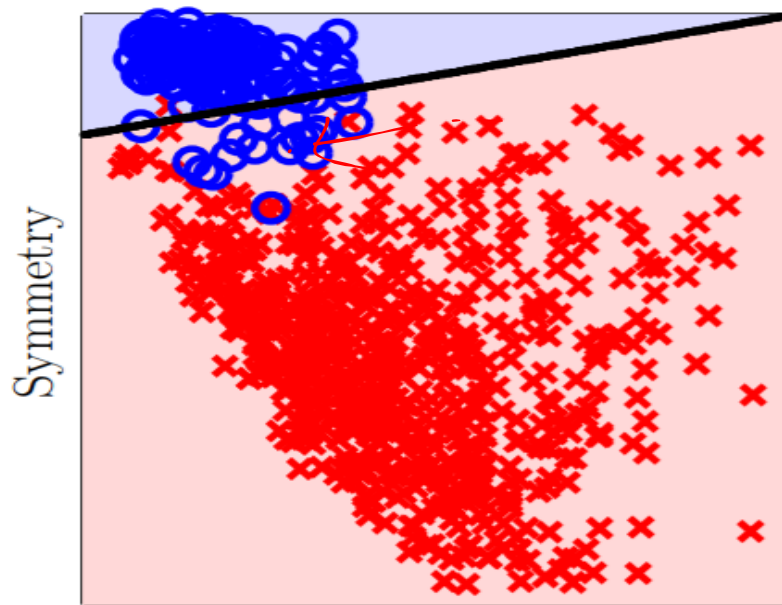
$$\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$$

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

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

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

# Digits Data

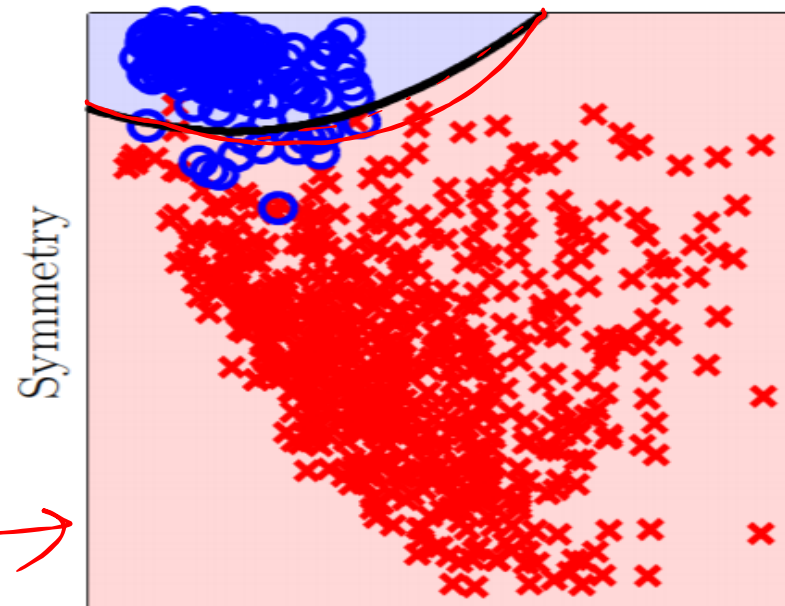


Average Intensity

**Linear model**

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

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



Average Intensity

**3rd order polynomial model**

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

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

10/0

# Humans Overfit (Superstitions)

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

Unfortunate events.

Overfitting lead to Opposite effect.

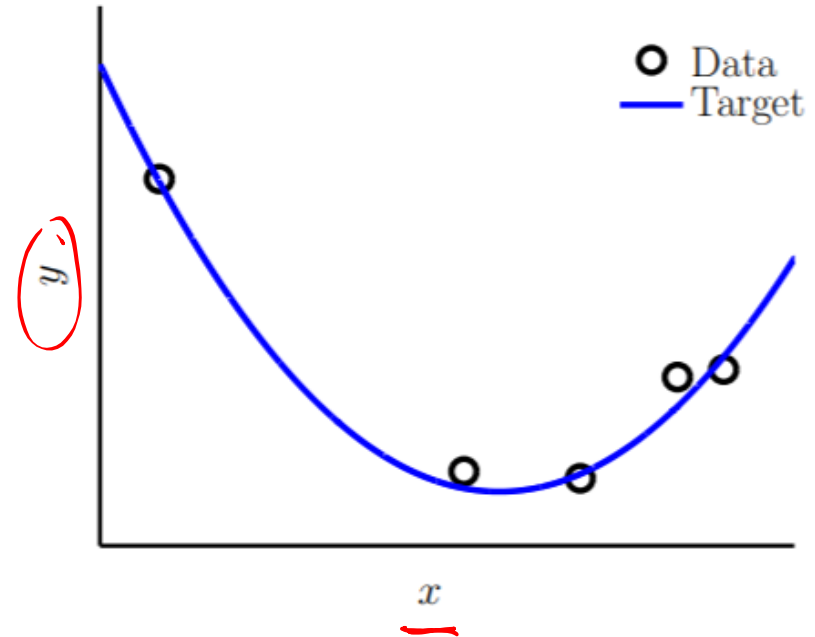
# Illustration of Overfitting

Quadratic  $f$

5 data points

A *little* noise (measurement error)

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



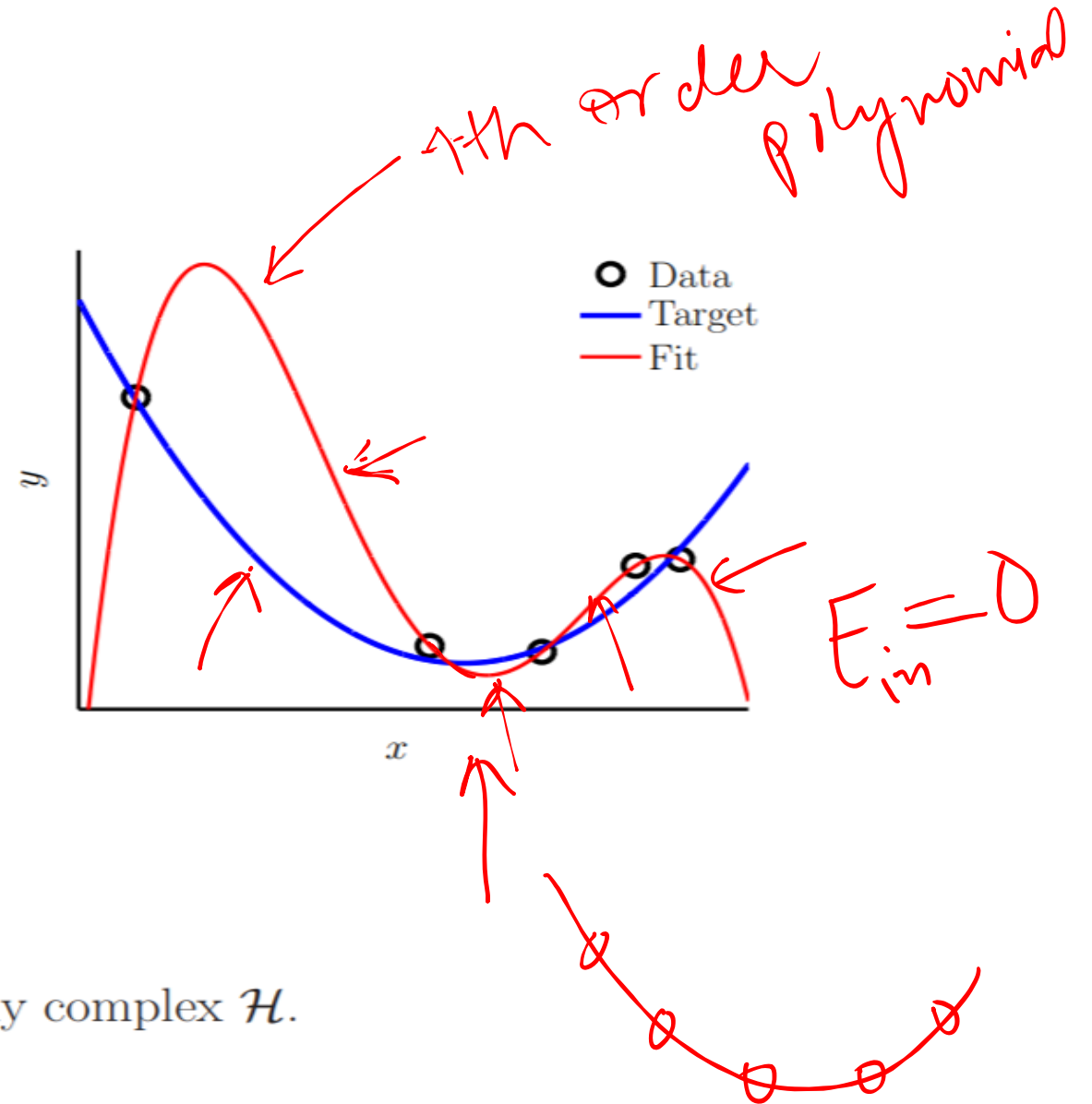
# Overfitting Example

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  $\mathcal{H}$ .

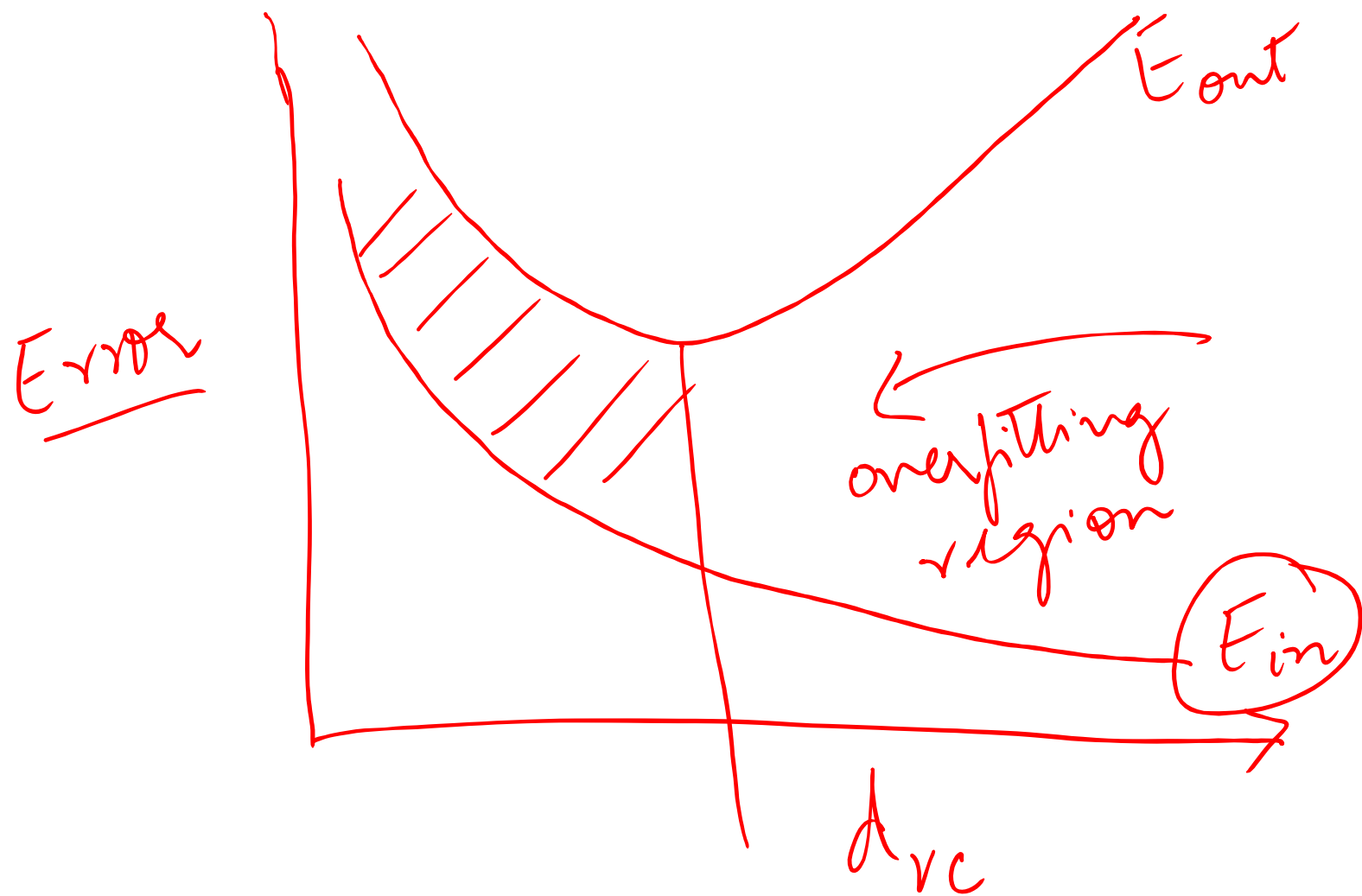
$$E_{in} \approx 0; E_{out} \gg 0$$

# Define Overfitting

- Fitting the data more than is warranted.

→ bad generalization  
→ process of picking a  $H$  with lower  $E_{in}$  resulting in a higher  $E_{out}$ .





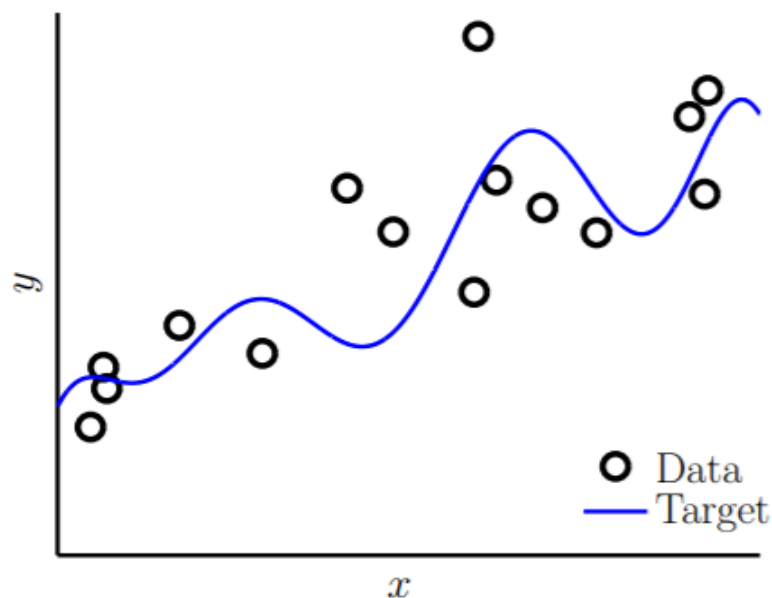
# Case Study

Choice of  $H$ -sets:  $H_2$  (2nd order polynomial)  
 $H_{10}$ : (10th order)

$$H_2 [1, n] \xrightarrow{\Phi} [1, x, x^2]$$

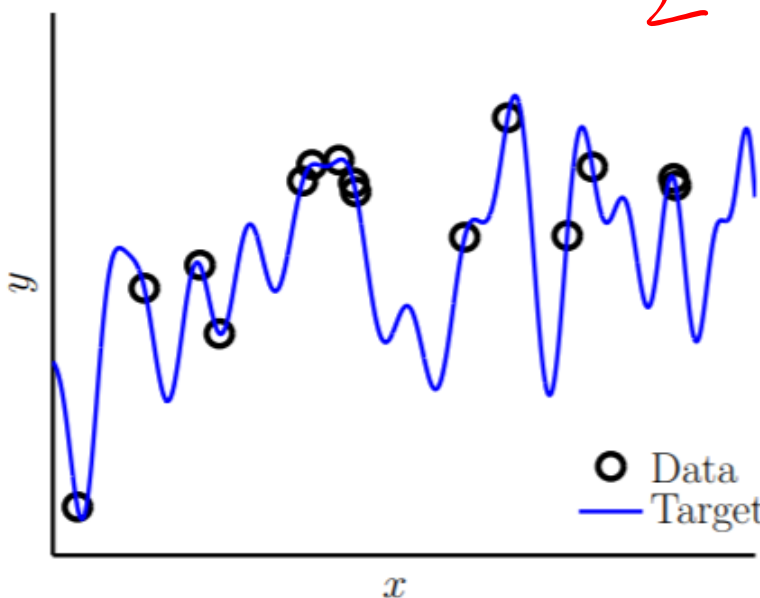
$$H_{10} [1, n] \xrightarrow{\Phi}$$

$$[1, x, x^2, \dots, x^{10}]$$



10th order  $f$  with noise.

(a)

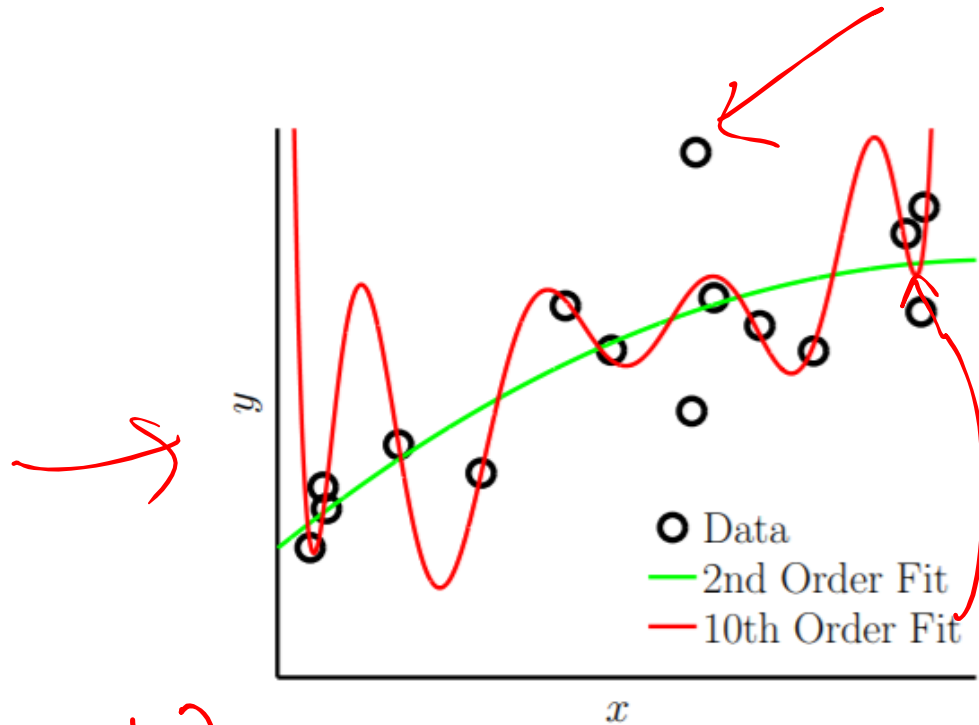


50th order  $f$  with no noise.

(b)

	f: 10th order with noise	f: 50th order no noise
<u>simple</u> $H_2$	noise X	No hope
<u>complex</u> $H_{10}$	✓ X	✓ X

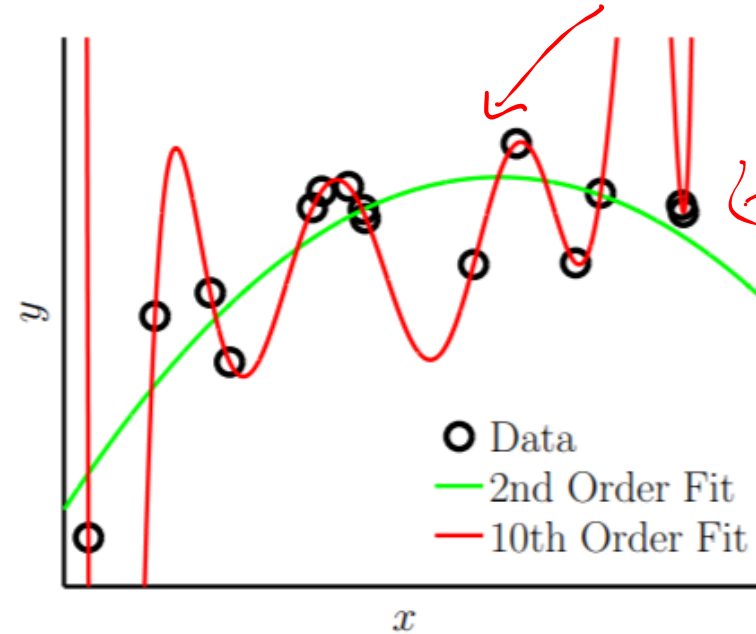
# 2<sup>nd</sup> order vs. 10<sup>th</sup> order polynomial



(a)

simple noisy target

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

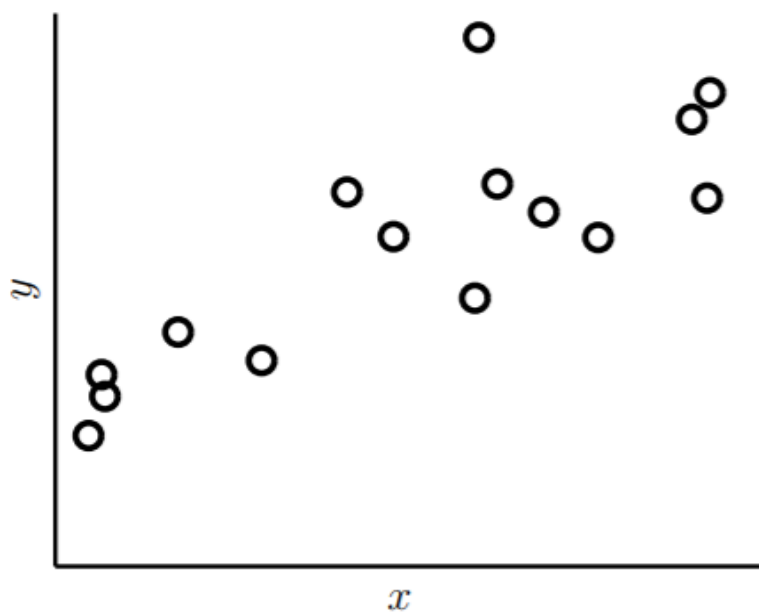


50th order

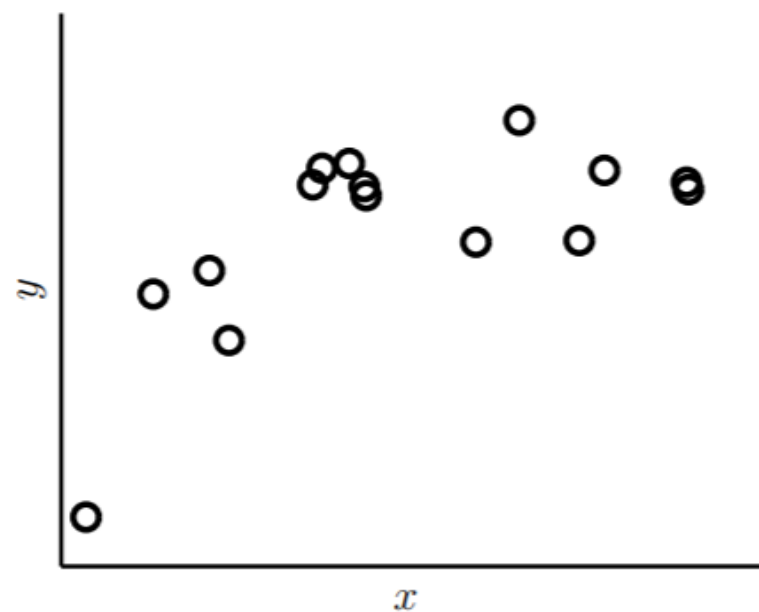
(b)

complex noiseless target

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



Simple  $f$  with noise.



Complex  $f$  with no noise.

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

# Measure Overfitting

•  $H_2 : E_{in}(H_2)$  ,  $H_{10} : E_{in}(H_{10})$

$$E_{in}(H_2) \geq E_{in}(H_{10})$$

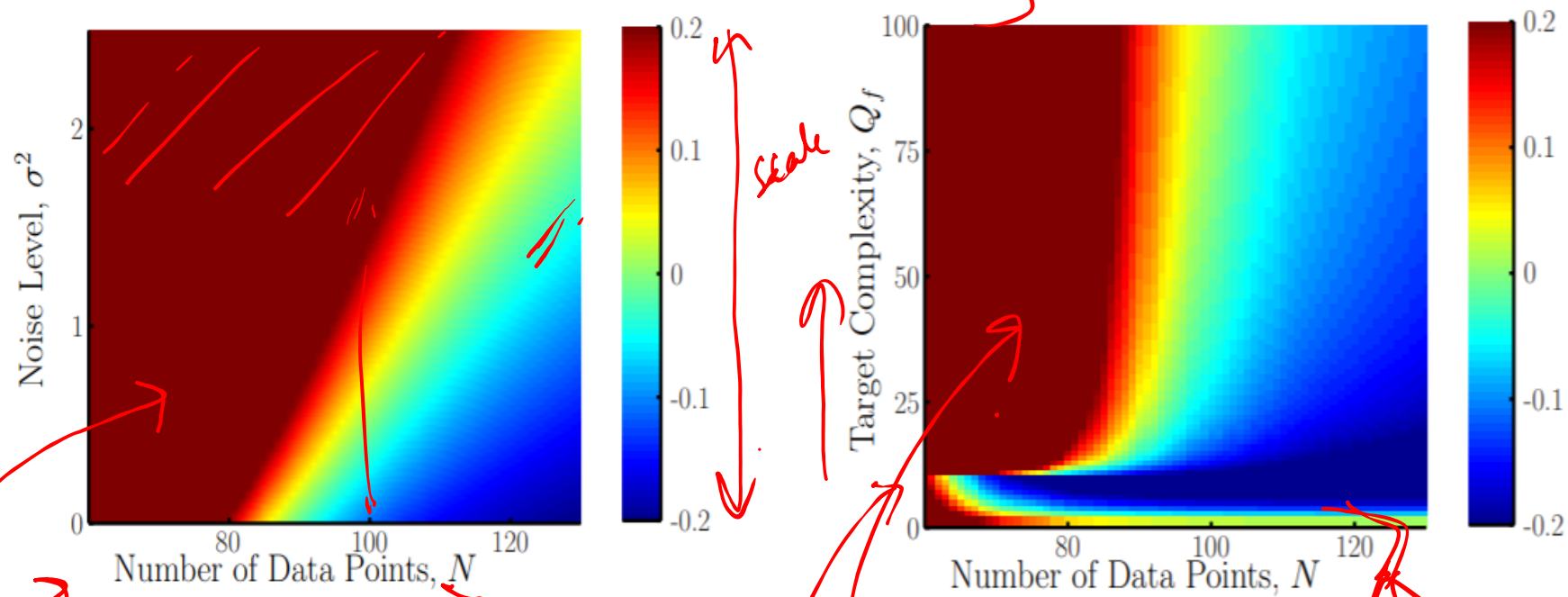
If  $E_{out}(H_{10}) > E_{out}(H_2)$   $\rightarrow g^{10}$

Overfit measure :  $E_{out}(H_{10}) - E_{out}(H_2)$   $\rightarrow g^2$

Controlled: i) Noise (measurement error)  $(\tilde{u})$   $N$

iii) complexity of  $f$   
(degree of poly.  $D$ )

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



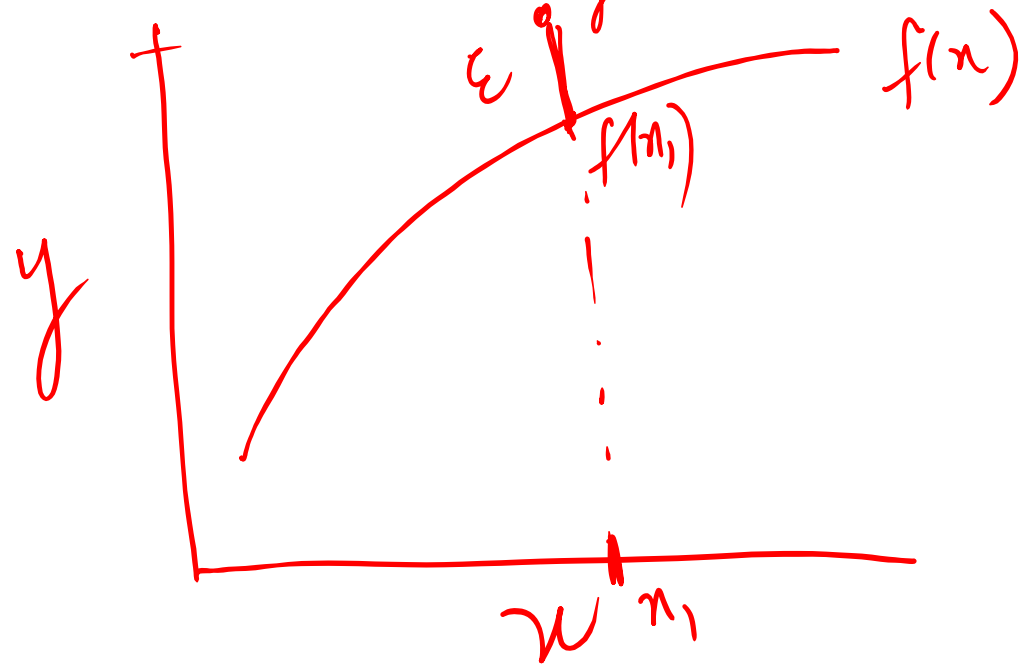
Number of data points $\uparrow$	Overfitting $\downarrow$
Noise $\uparrow$	Overfitting $\uparrow$
Target complexity $\uparrow$	Overfitting $\uparrow$

# Noise

- That part of  $y$  which cannot be modelled  
↳ observed value.

↳ 2 Sources:  
i) Stochastic Noise.  
(measurement error)

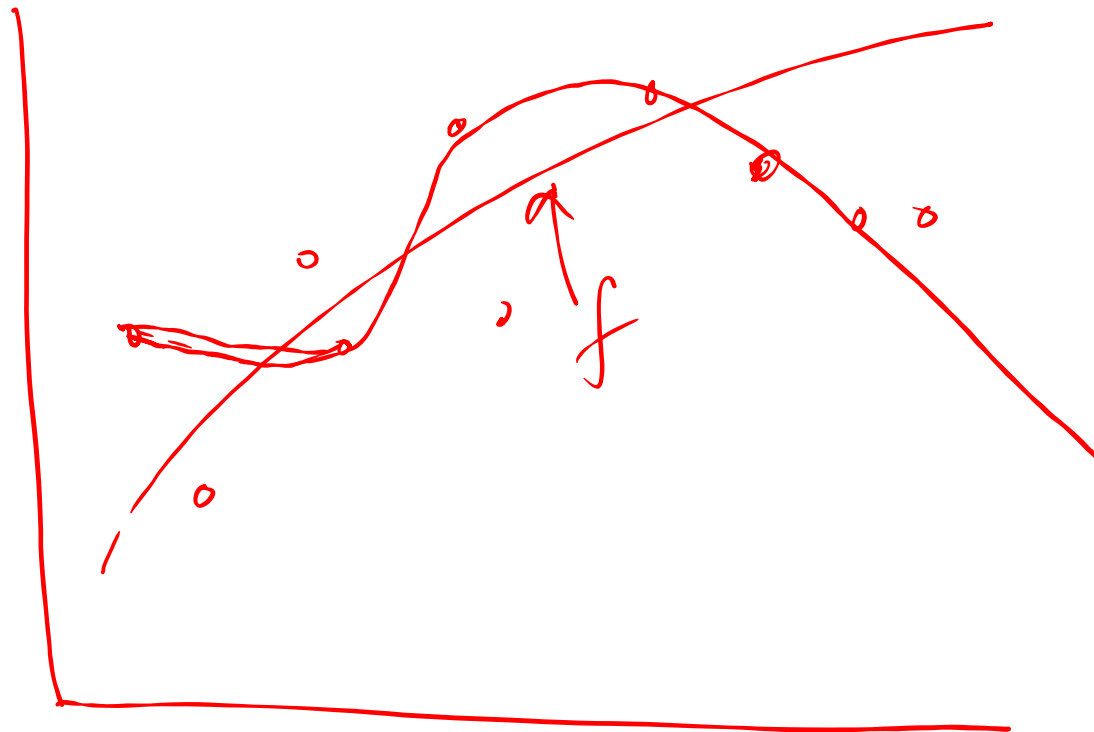
$$y = f(x) + \epsilon$$
$$y = f(x)$$





Data with noise  
→ higher chances of  
my model  
going astray.

Noise → cannot be  
modelled



# Stochastic Noise

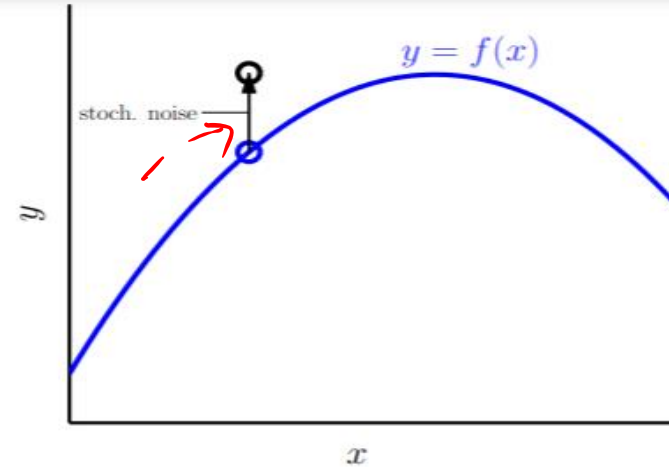
We would like to learn from ○:

$$y_n = f(x_n)$$

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.



100 pictures

# Deterministic Noise

→ fixed N-set

that part of the  $(f)$  that  
cannot be modelled.

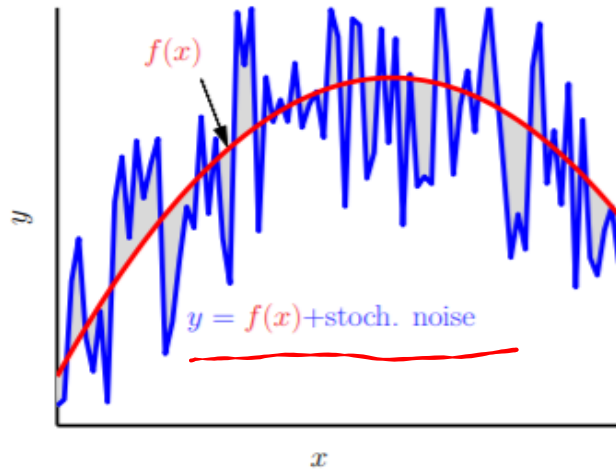


Stochastic & Determinist ~~effect~~  $\rightarrow$  Overfit.  
Regenerate  $n$  values  $\begin{cases} S \rightarrow \text{different } y \text{ values} \\ D \rightarrow \text{same } y \text{ values.} \end{cases}$

ii)  $S \xrightarrow{\text{deficiency}} \mathcal{N}$

$D \xrightarrow[\text{a complex}]{\text{choose}} \mathcal{N} \} \uparrow \text{Trade-off.}$

## Stochastic Noise ✓



source: random measurement errors

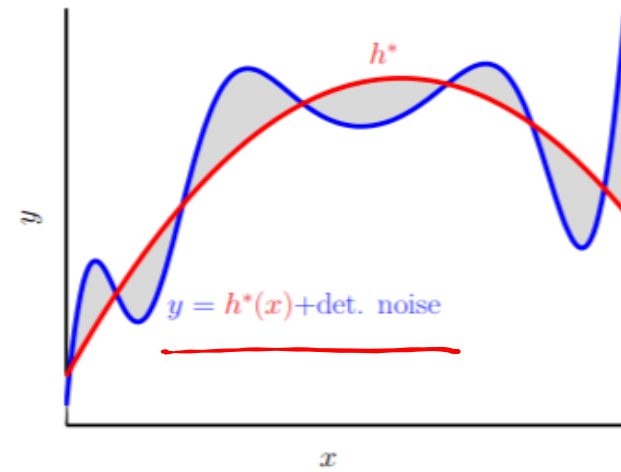
✓ re-measure  $y_n$

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

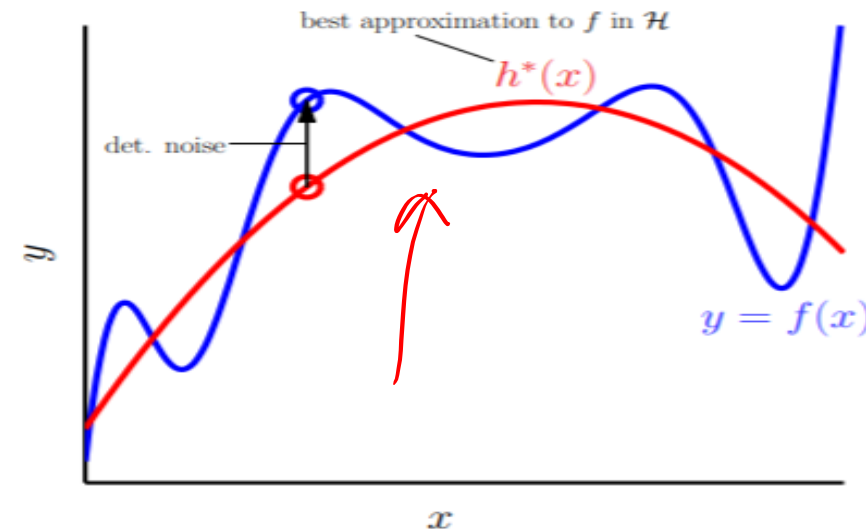
# Deterministic Noise

We would like to learn from  $\circ$ :

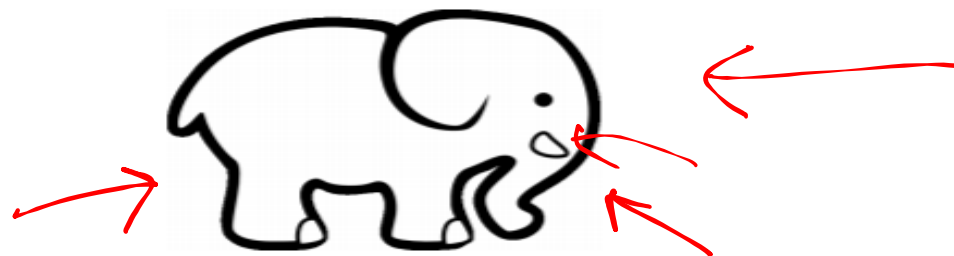
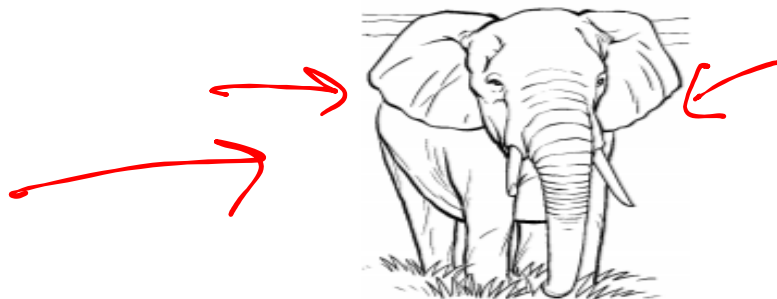
$$y_n = h^*(x_n)$$

Unfortunately, we only observe  $\circ$ :

$$\begin{aligned} y_n &= f(x_n) \\ &= h^*(x_n) + \text{'deterministic noise'} \\ &\quad \uparrow \\ &\quad \mathcal{H} \text{ cannot model this} \end{aligned}$$



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



# Bias Variance Analysis and Noise

- $$\underline{E_D[E_{out}]} = \underbrace{\text{Bias}(x)}_{(g(x) - f(x))^2} + \frac{\text{Variance}(x)}{\text{var}(g(x))}$$

$$\underline{E_{out}(x)} = (g(x) - y)^2, \quad y = f(x) + \epsilon$$
$$= (g(x) - f(x) - \epsilon)^2$$

$$E_D[E_{out}(x)] = \underbrace{(g(x) - f(x))^2} - 2 \cancel{\epsilon g(x) f(x)} + \underbrace{\epsilon^2}$$

$$\underline{E[E_{out}^{(n)}]} = \underbrace{E(g(x) - f(x))^2}_{\text{Bias + Variance}} + \underbrace{E(\epsilon^2)}_{\sigma^2}$$

$$= \text{Bias} + \text{Variance} + \sigma^2$$

(n & N) N ↑

$$E[E_{out}] = \underbrace{\text{Bias}}_{\substack{\swarrow \\ \text{deterministic} \\ \text{noise} \\ \checkmark}} + \underbrace{\text{Variance}}_{\substack{\text{indirect} \\ \text{impact of} \\ \text{noise}}} + \sigma^2 \rightarrow \text{stochastic} \\ \text{noise} \\ \checkmark$$

Model  $\rightarrow$  lower bias, lower variance



## Summarize

Overfitting → known problem.  
Noise → cannot be modelled, always there.  
→ Can we do something about this?  
→ Yes → Regularization.

Validation

# Summary

-

Thanks!