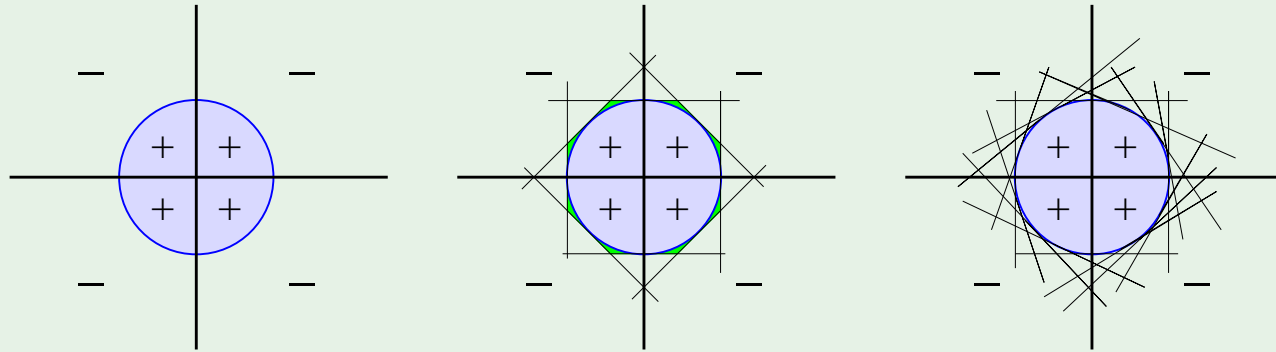


# Review of Lecture 10

- Multilayer perceptrons

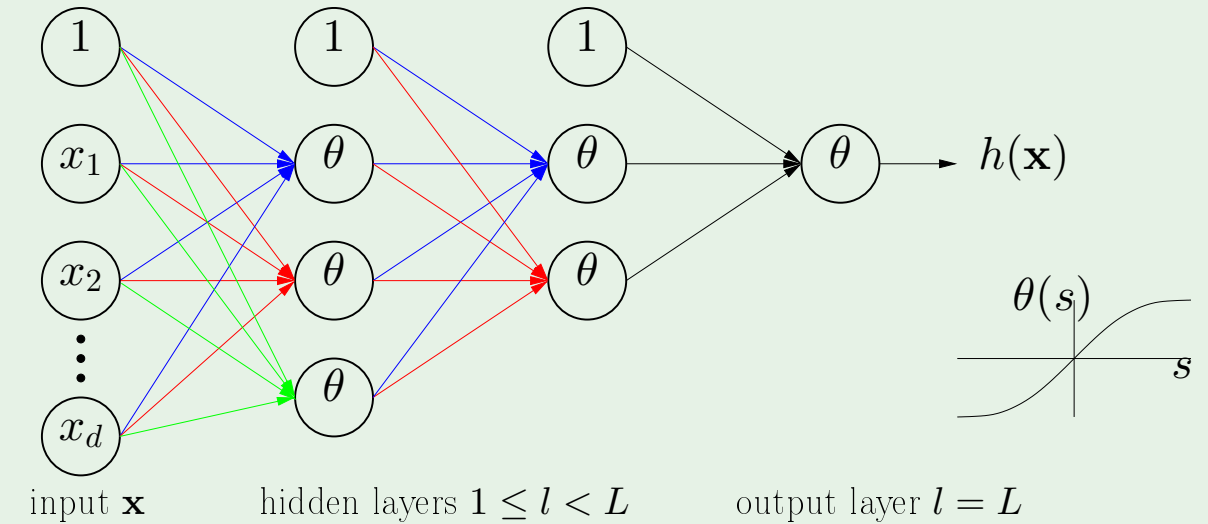


Logical combinations of perceptrons

- Neural networks

$$x_j^{(l)} = \theta \left( \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)} \right)$$

where  $\theta(s) = \tanh(s)$



- Backpropagation

$$\Delta w_{ij}^{(l)} = -\eta x_i^{(l-1)} \delta_j^{(l)}$$

where

$$\delta_i^{(l-1)} = (1 - (x_i^{(l-1)})^2) \sum_{j=1}^{d^{(l)}} w_{ij}^{(l)} \delta_j^{(l)}$$

Most important aspect of tanh is that it is smooth/twice differentiable, so the dependency of the error in the output on w will be a well-behaved function and we can apply gradient descent

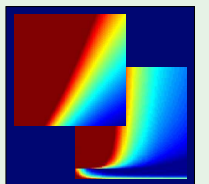
# Learning From Data

Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 11: Overfitting



Sponsored by Caltech's Provost Office, E&AS Division, and IST • Tuesday, May 8, 2012



# Outline

- What is overfitting?
- The role of noise
- Deterministic noise
- Dealing with overfitting

# Illustration of overfitting

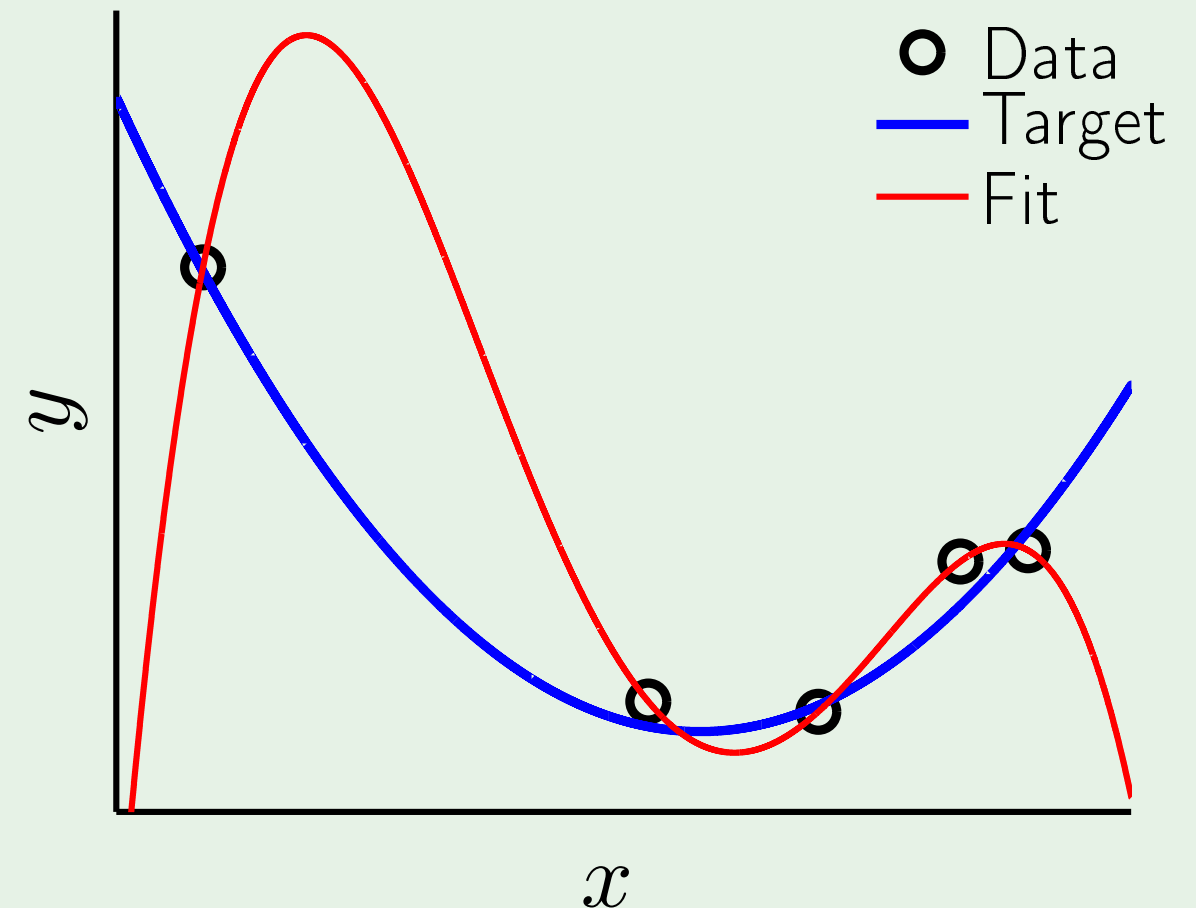
Simple target function

5 data points- **noisy**

4th-order polynomial fit

$$E_{\text{in}} = 0, \quad E_{\text{out}} \text{ is huge}$$

Overfitting is a comparative term, so one situation must be worse than another - you went further than you should have. There is a distinction between overfitting and bad generalization: since using a 3rd order polynomial we would not in general get zero  $E_{\text{in}}$  but we would get a better  $E_{\text{out}}$ , so the bad generalization was a result of using 4th order rather than 3rd order, hence overfitting.



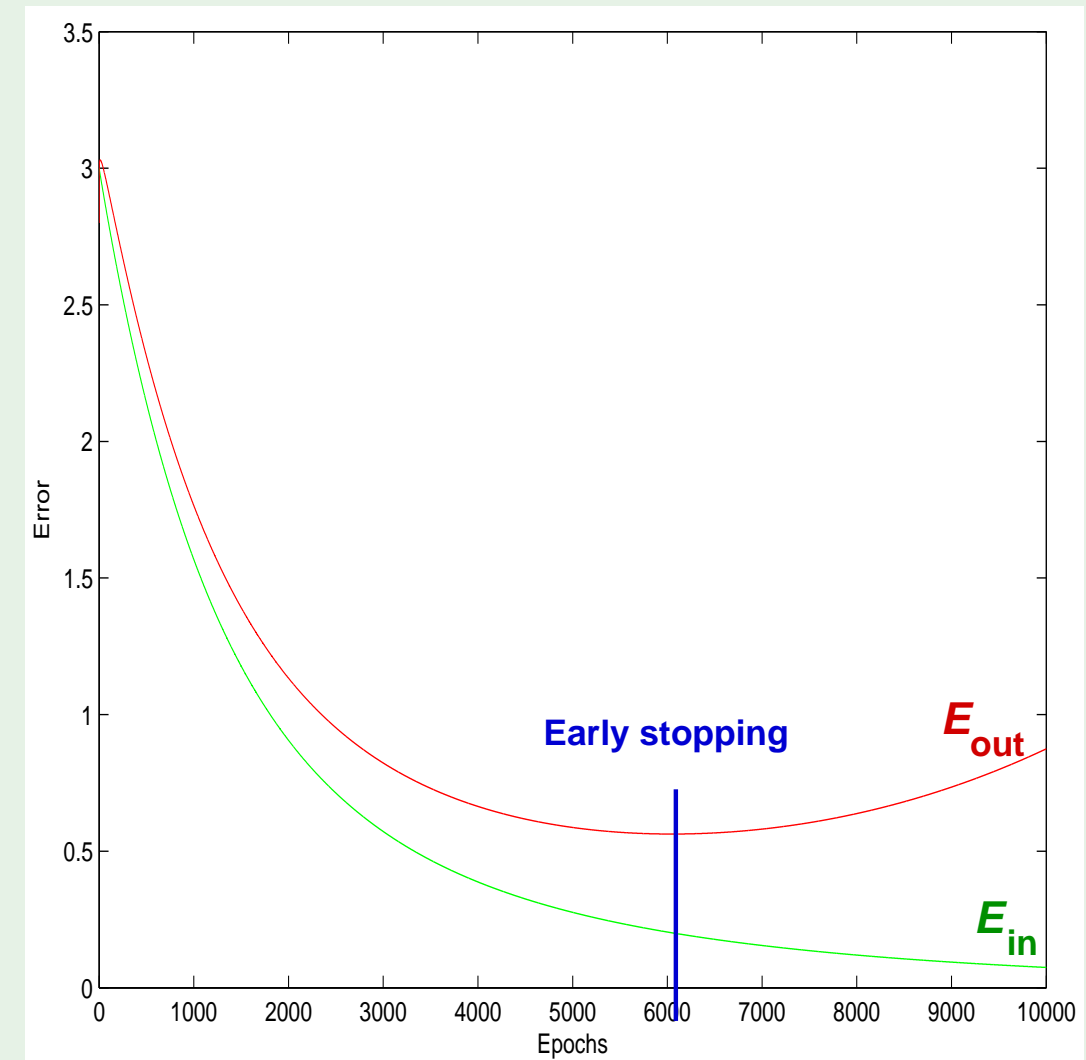
Note that this example shows overfitting between two models. However overfitting in neural networks is within the same model.

# Overfitting versus bad generalization

Neural network fitting noisy data

Overfitting: occurs when  $E_{\text{in}} \downarrow$  but  $E_{\text{out}} \uparrow$

$E_{\text{in}}$  and  $E_{\text{out}}$  are similar values at the start since we only have one hypothesis which does not depend on the dataset, the randomly initialized weights, so not a surprise they are the same value. As the epochs increase and we explore more of the weight space, we get the associated benefit/harm from this increase gradually - the effective VC dimension (if we can define that) grows as we explore the space until we have (potentially) explored the whole space, ending as the total number of free parameters of the model.



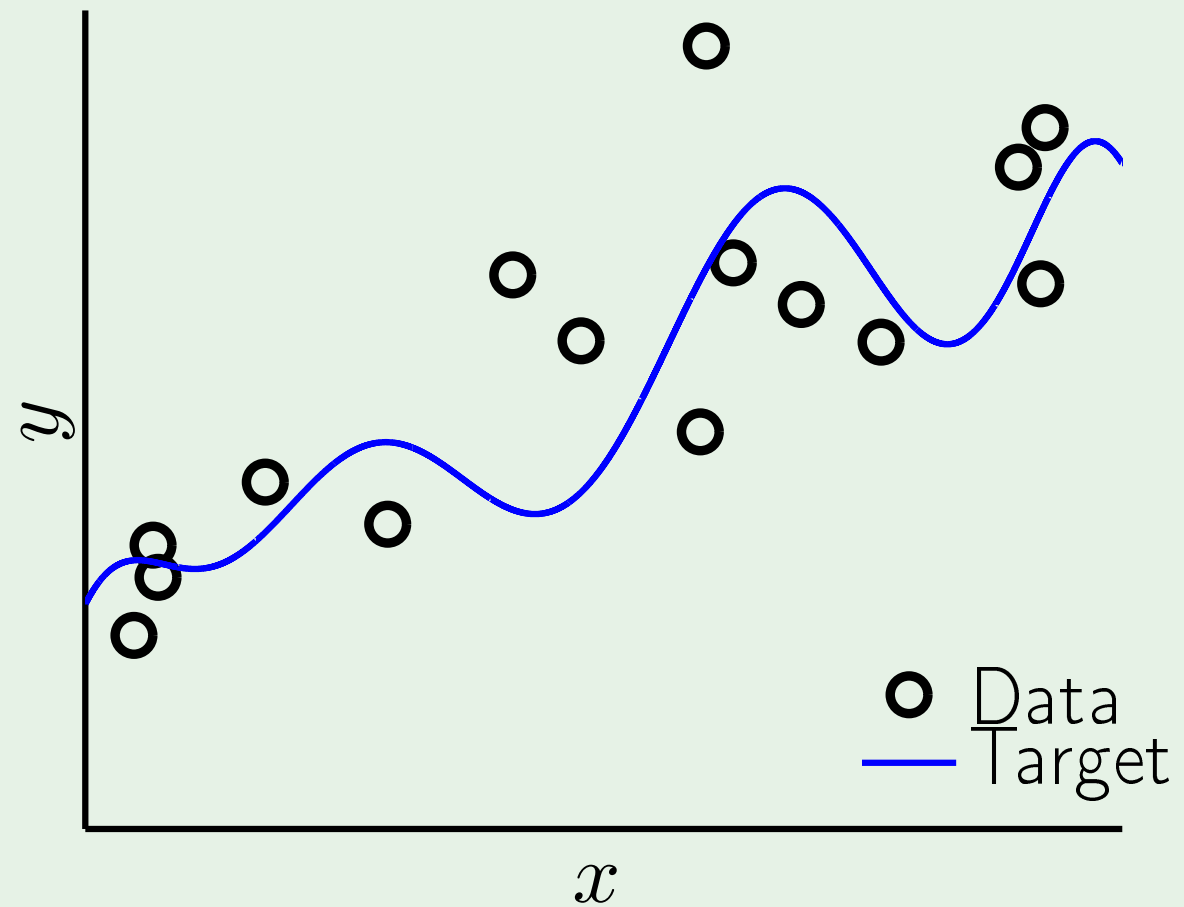
# The culprit

**Overfitting:** “fitting the data more than is warranted”

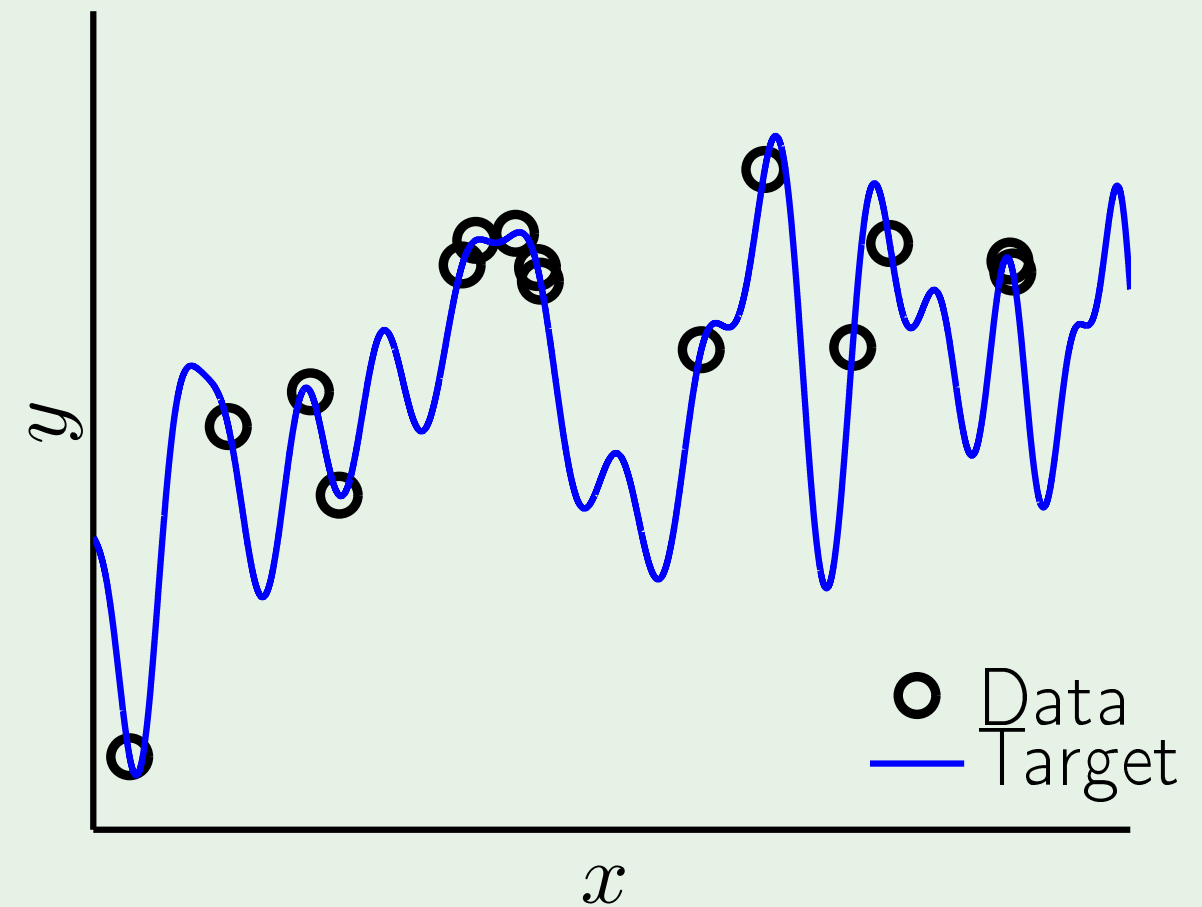
**Culprit:** fitting the noise - harmful

# Case study

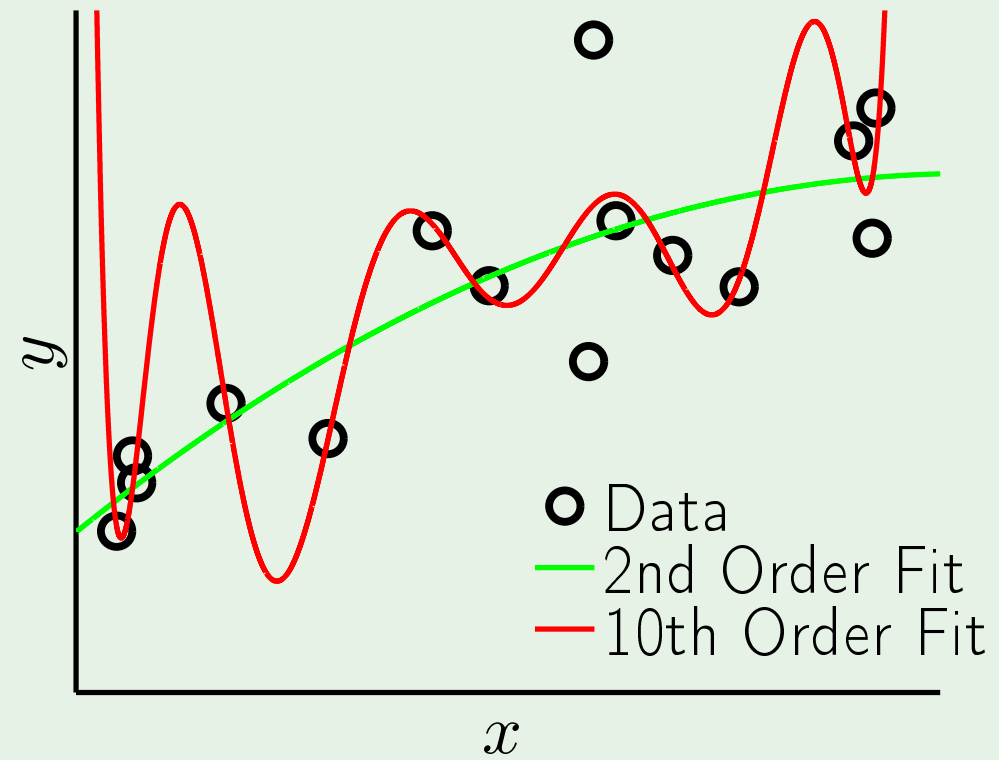
10th-order target + noise



50th-order target

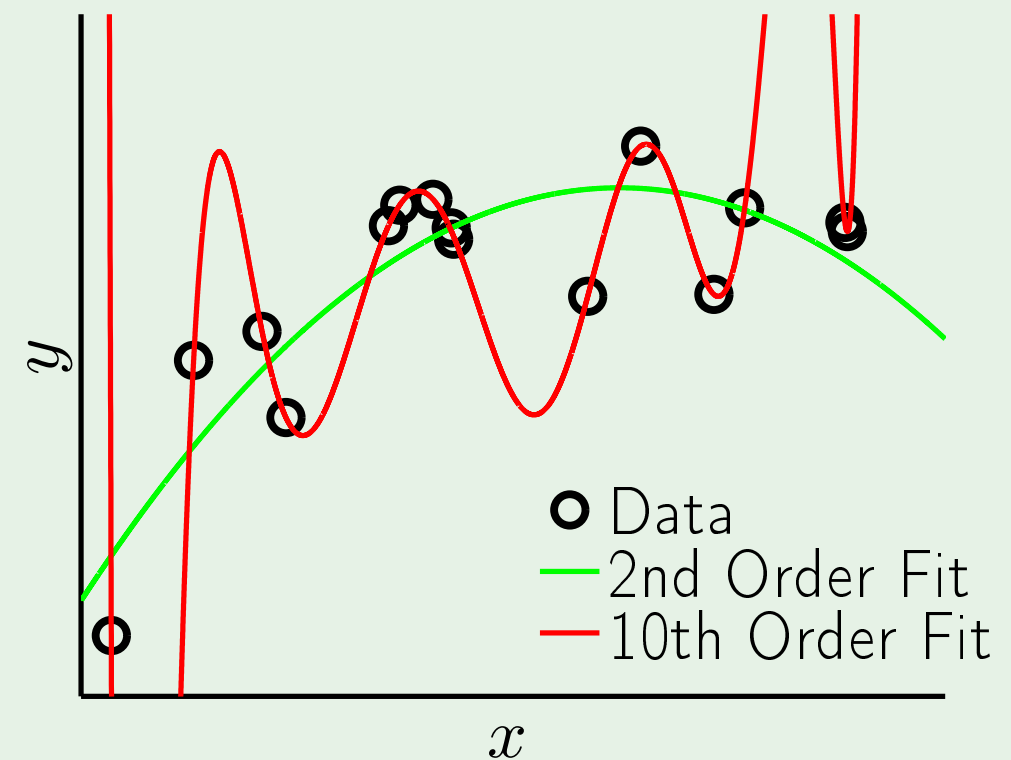


## Two fits for each target



Noisy low-order target

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



Noiseless high-order target

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



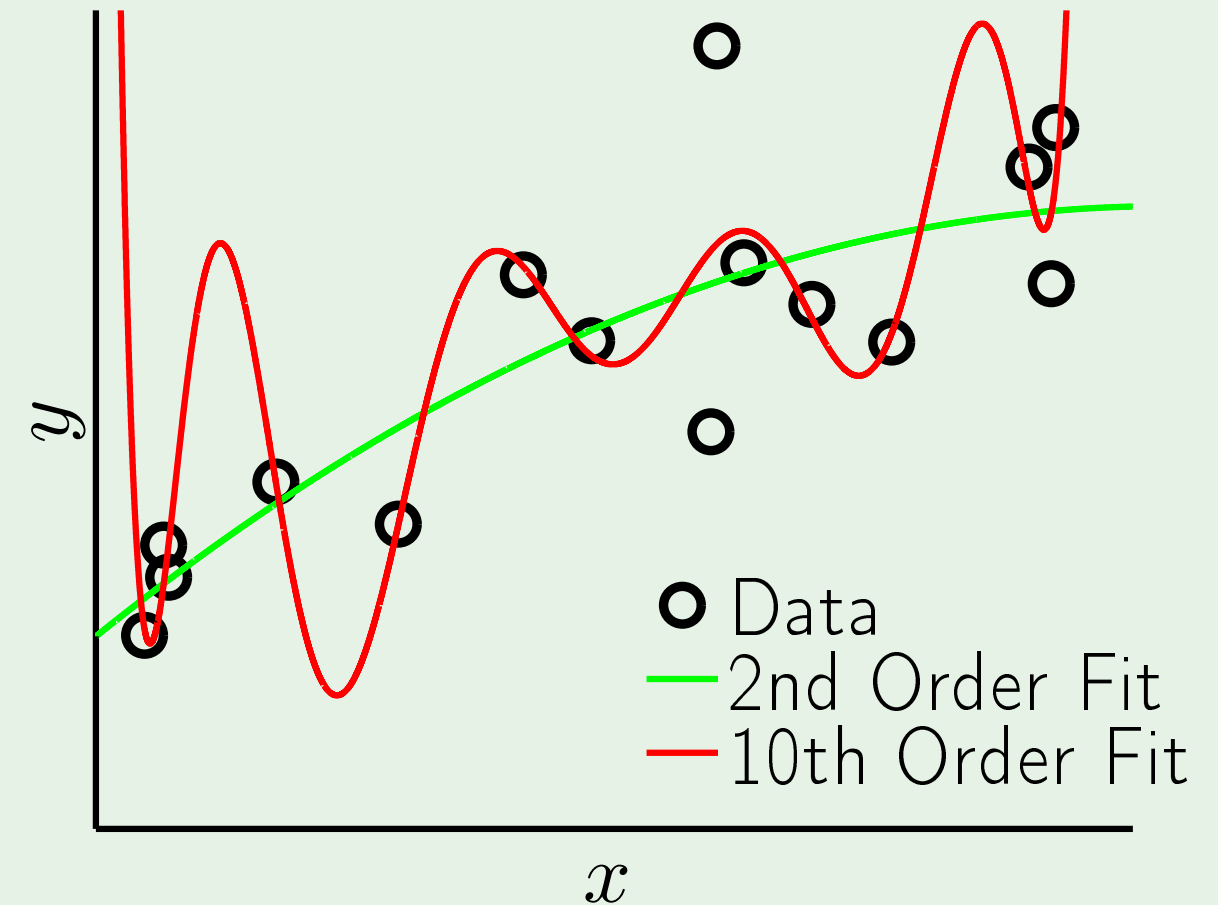
# An irony of two learners

Two learners  $O$  and  $R$

They know the target is 10th order

$O$  chooses  $\mathcal{H}_{10}$        $R$  chooses  $\mathcal{H}_2$

Better to match complexity of the model to the amount of data we have - generalization issues depend on the size and the quality of the dataset.

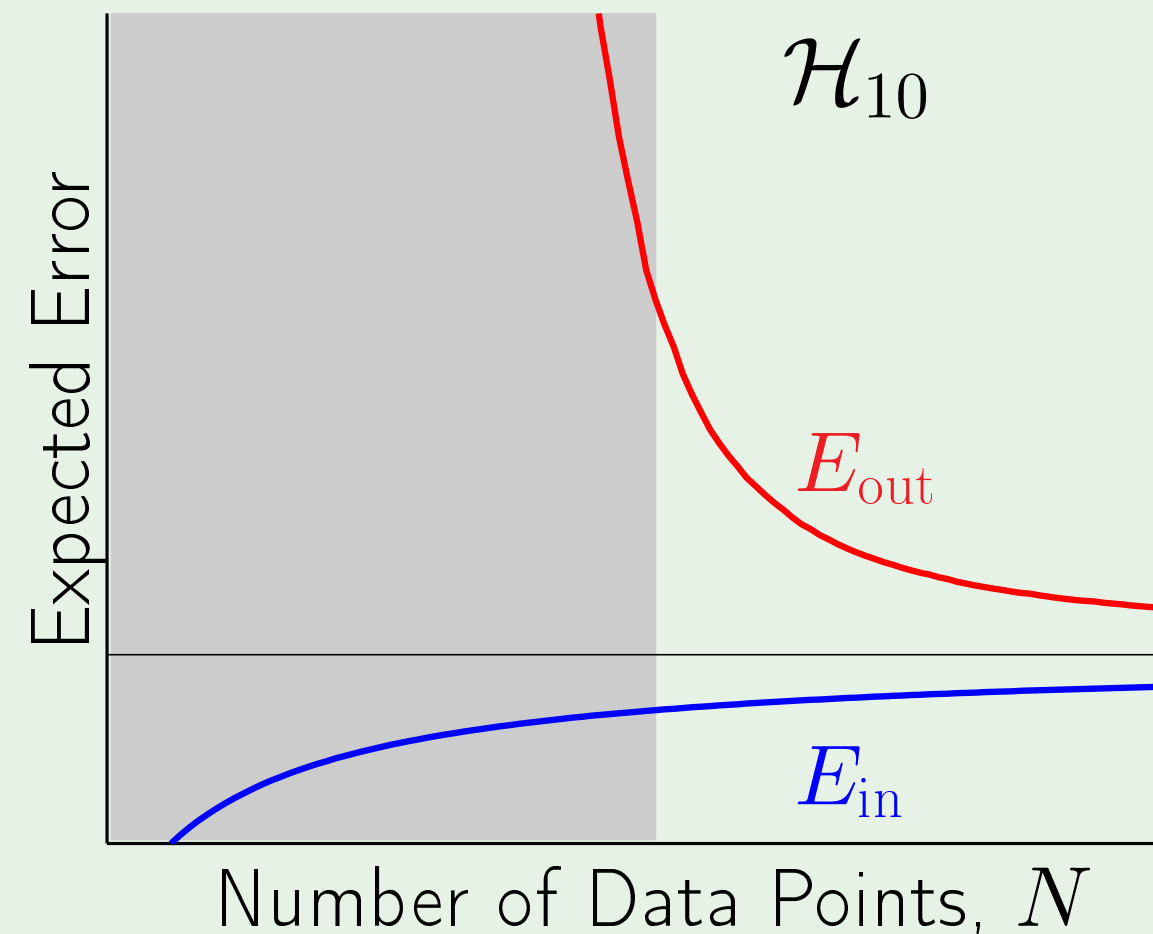
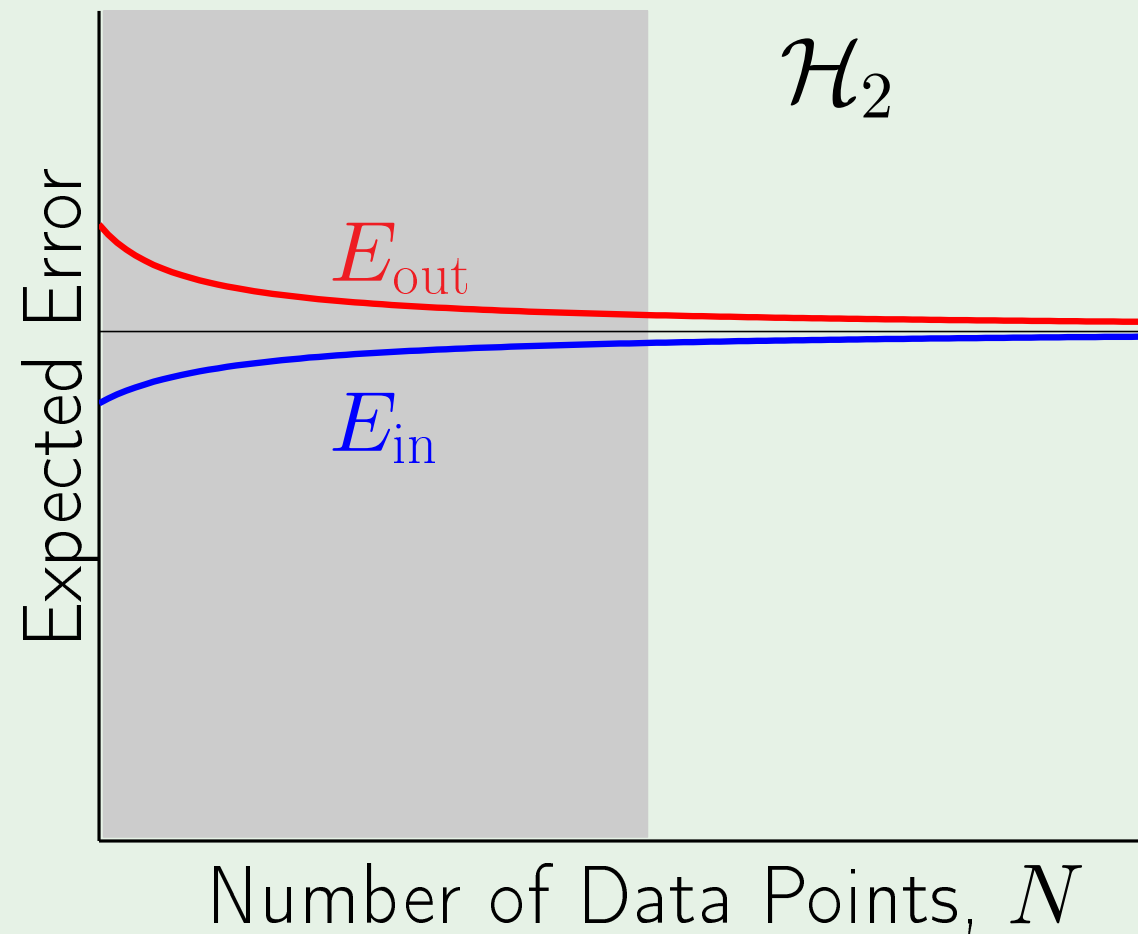


Learning a 10th-order target

## We have seen this case

Remember learning curves?

The baseline/inevitable error comes from both the limitations of the model and the noise in the data. The grey area is where overfitting is happening: while  $E_{in}$  for  $H_{10}$  is smaller (even if it is smaller for all  $N$ ),  $E_{out}$  for  $H_{10}$  is larger than for  $H_2$ .



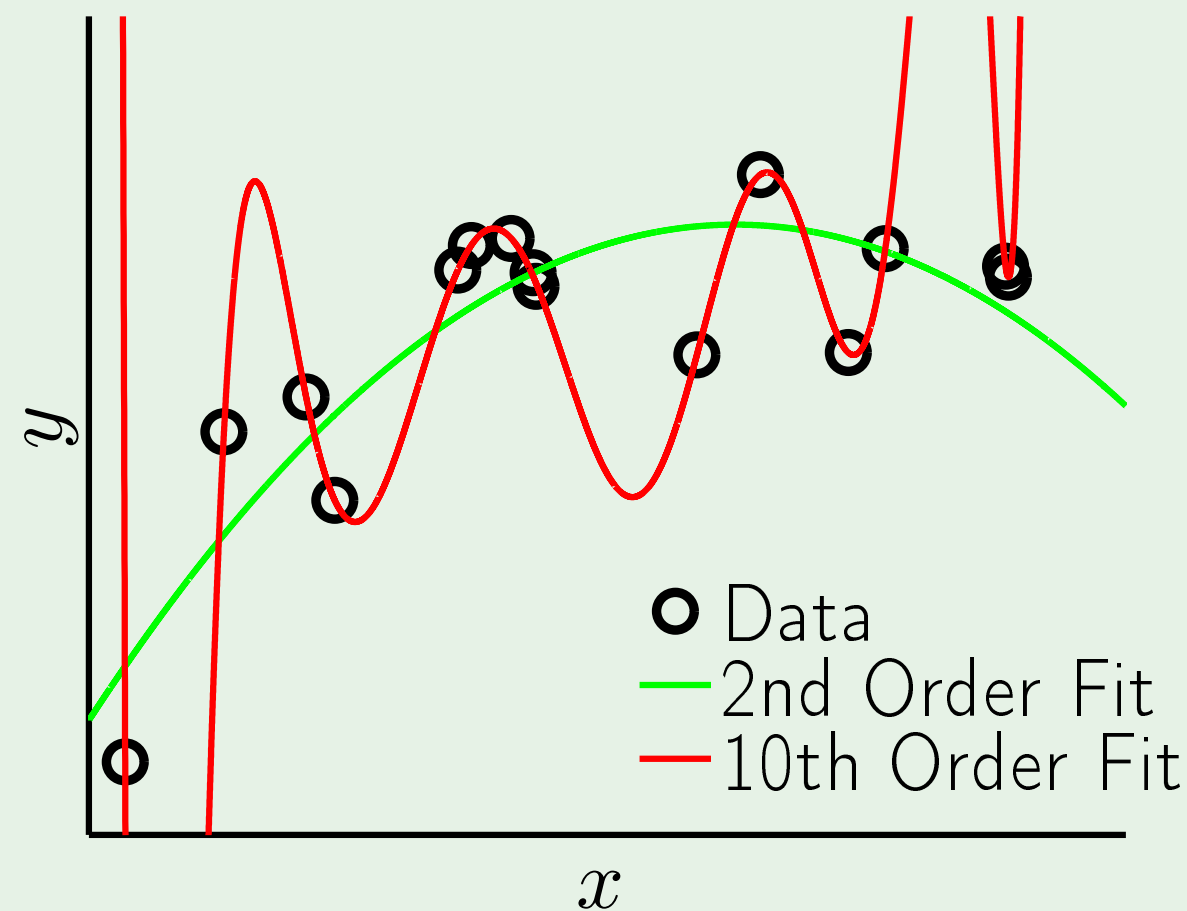
## Even without noise

The two learners  $\mathcal{H}_{10}$  and  $\mathcal{H}_2$

They know there is no noise

Is there really no noise?

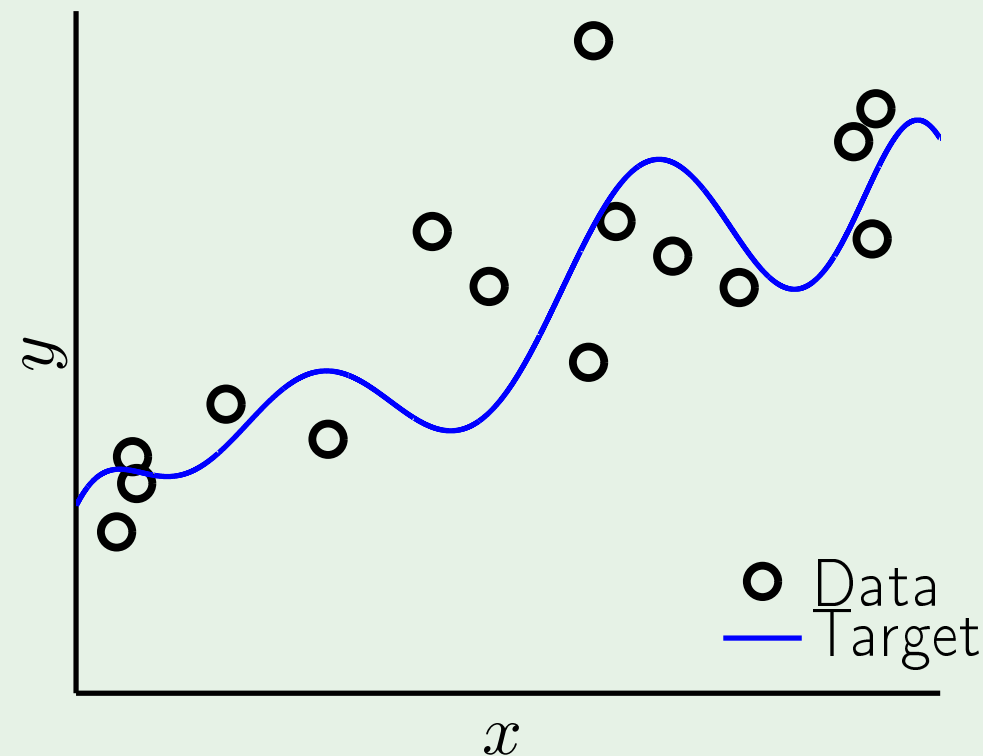
It seems that overfitting occurs both when we have noise in the data and when the data appears noiseless, so it seems target complexity may be related - see next slide.



Learning a 50th-order target

# A detailed experiment

## Impact of noise level and target complexity



We do not generate the alpha coefficients randomly as doing that gives us non-interesting polynomials that are generally dominated by a single power of  $x$ . Instead of random alpha, we use normalized Legendre polynomials (i.e. the choice of the coefficients is such that two Legendre polynomials of different order are orthogonal) - like harmonics in a sinusoidal expansion, the inner product of each order Legendre polynomial are orthogonal. If we have a combination of Legendre with random coefficients, we can generate interesting shapes. The alphas now just become to sum of the coefficients for each power of  $x$ . We are now just generating the alphas in an elaborate way so that we get interesting targets.

$$y = f(x) + \underbrace{\epsilon(x)}_{\sigma^2} = \sum_{\substack{q=0 \\ \text{normalized}}}^{Q_f} \alpha_q x^q + \epsilon(x)$$

noise level:  $\sigma^2$  - with expectation value (of the error function) as zero

target complexity:  $Q_f$  - describes the complexity of  $f$

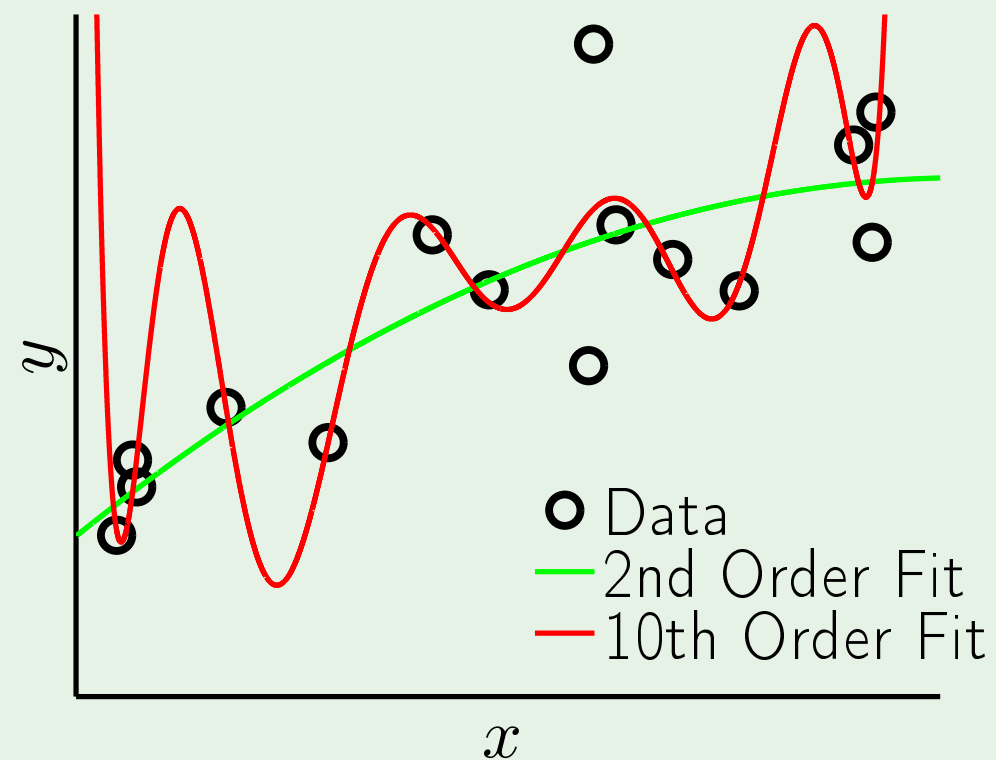
data set size:  $N$

# The overfit measure

We fit the data set  $(x_1, y_1), \dots, (x_N, y_N)$  using our two models:

$\mathcal{H}_2$ : 2nd-order polynomials

$\mathcal{H}_{10}$ : 10th-order polynomials



Compare out-of-sample errors of

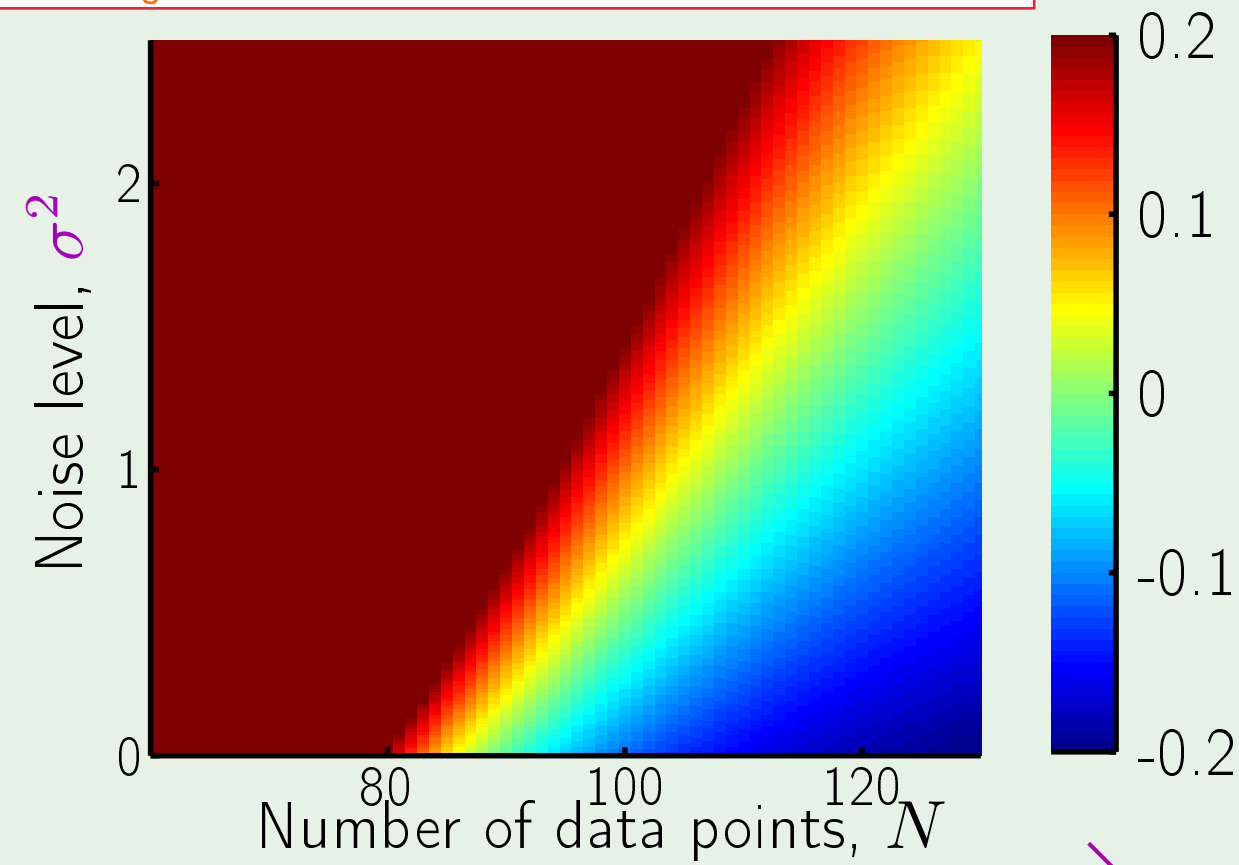
$g_2 \in \mathcal{H}_2$  and  $g_{10} \in \mathcal{H}_{10}$

**overfit measure:**  $E_{\text{out}}(g_{10}) - E_{\text{out}}(g_2)$

Large positive value indicates a lot of overfitting, small positive means less overfitting, negative indicates no overfitting

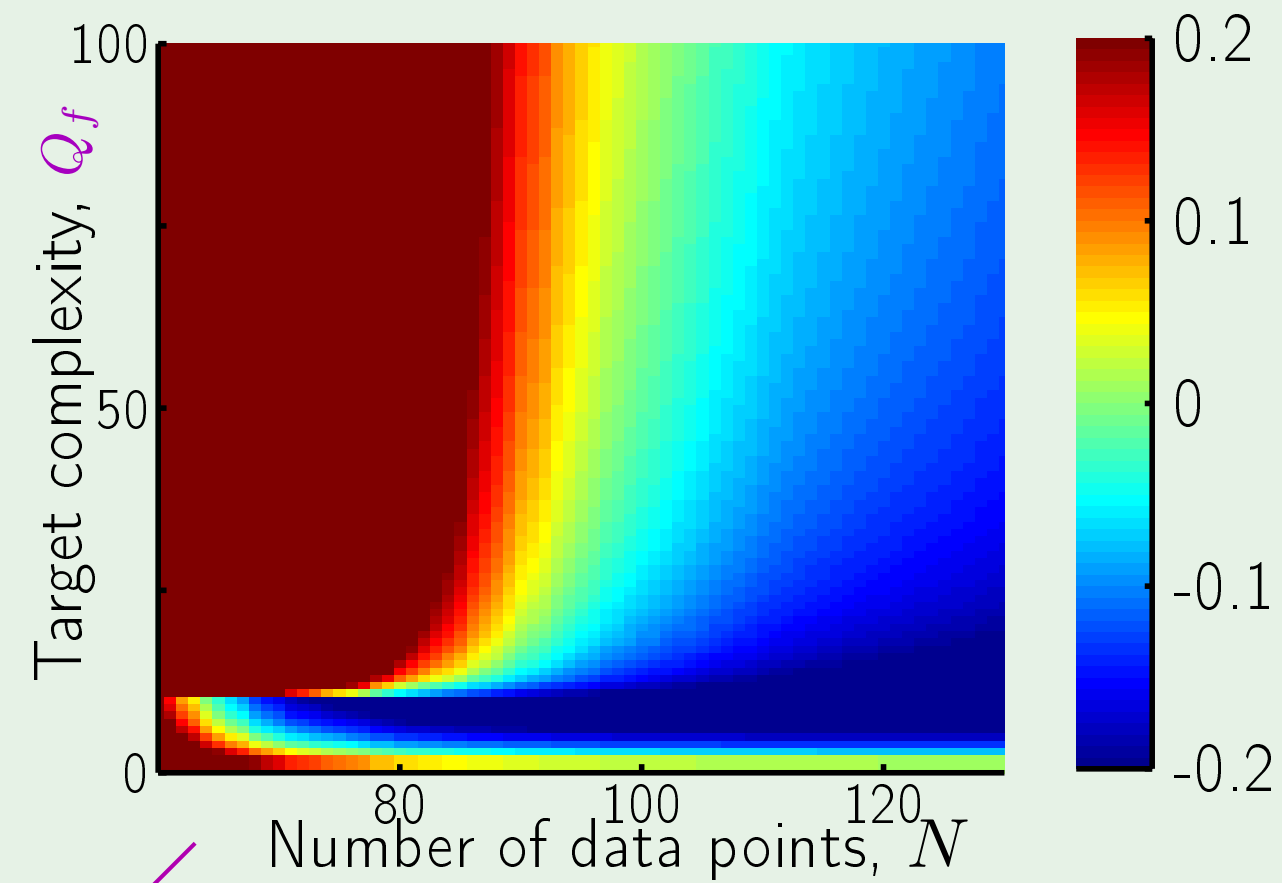
Here we do the experiment over ten million iterations (for all kinds of parameters): generate  $f$ , generate a dataset, fit using both hypothesis sets, evaluate the overfit measure. Here the colourbar indicates the intensity of overfit (the overfit measure) - depending on the number of data points  $N$  and the level of the noise. Red indicates more overfitting, and blue indicates no overfitting.

# The results



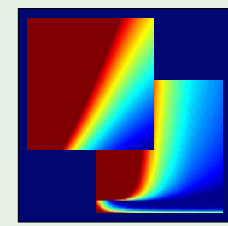
Impact of  $\sigma^2$

In this graph,  $Q_f$  is fixed at 20

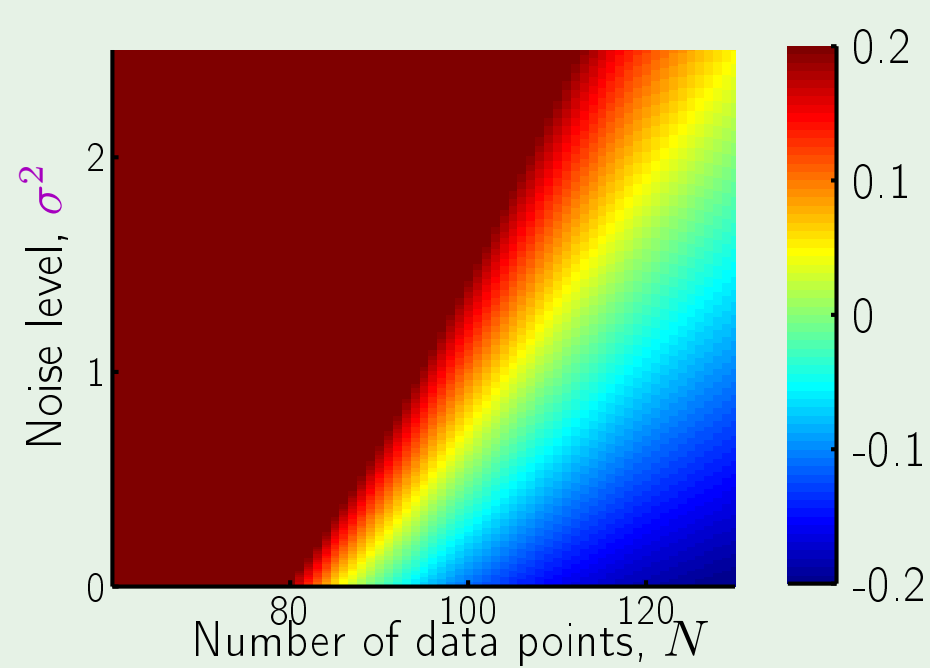


Impact of  $Q_f$

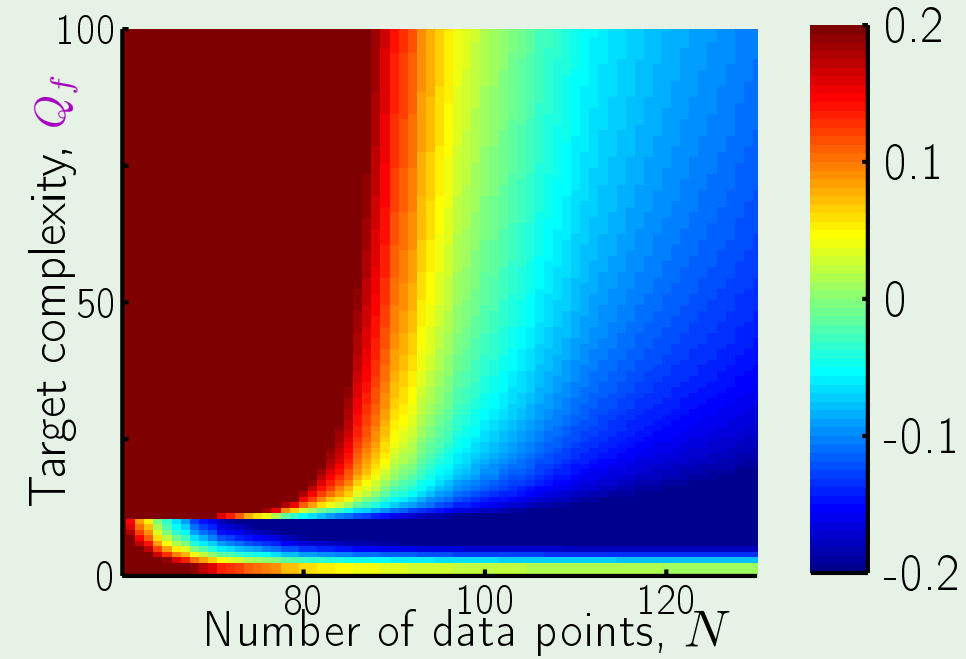
In this graph, the noise level is fixed at  $\sigma^2 = 0.1$



# Impact of “noise”



Stochastic noise



Deterministic noise

number of data points	↑	Overfitting	↓
stochastic noise	↑	Overfitting	↑
deterministic noise	↑	Overfitting	↑

# Outline

- What is overfitting?
- The role of noise
- Deterministic noise
- Dealing with overfitting



# Definition of deterministic noise

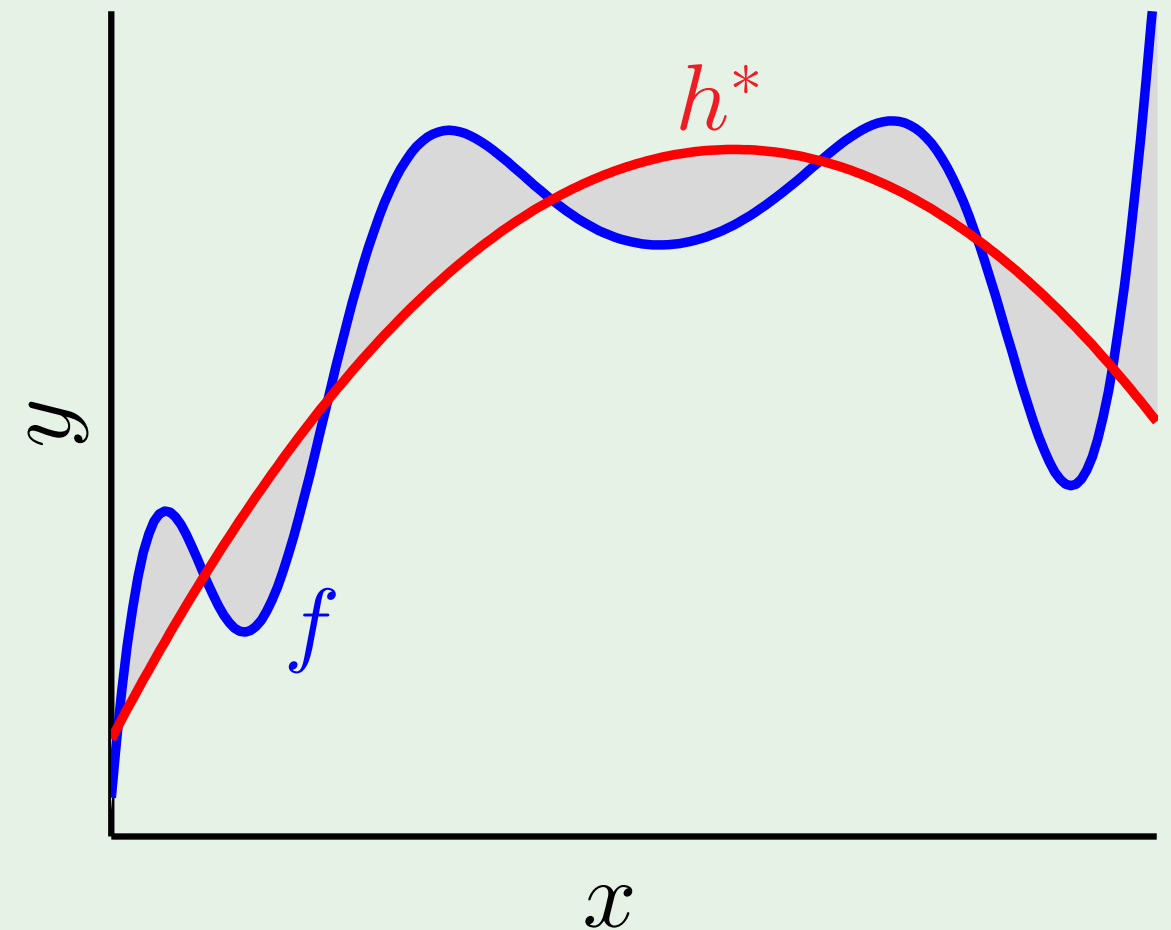
The part of  $f$  that  $\mathcal{H}$  cannot capture:  $f(\mathbf{x}) - h^*(\mathbf{x})$

So deterministic noise "looks like noise" to some hypothesis sets but does not look like noise to others and experimentally this affects overfitting

Why "noise"?

Main differences with stochastic noise:

1. depends on  $\mathcal{H}$
2. fixed for a given  $\mathbf{x}$

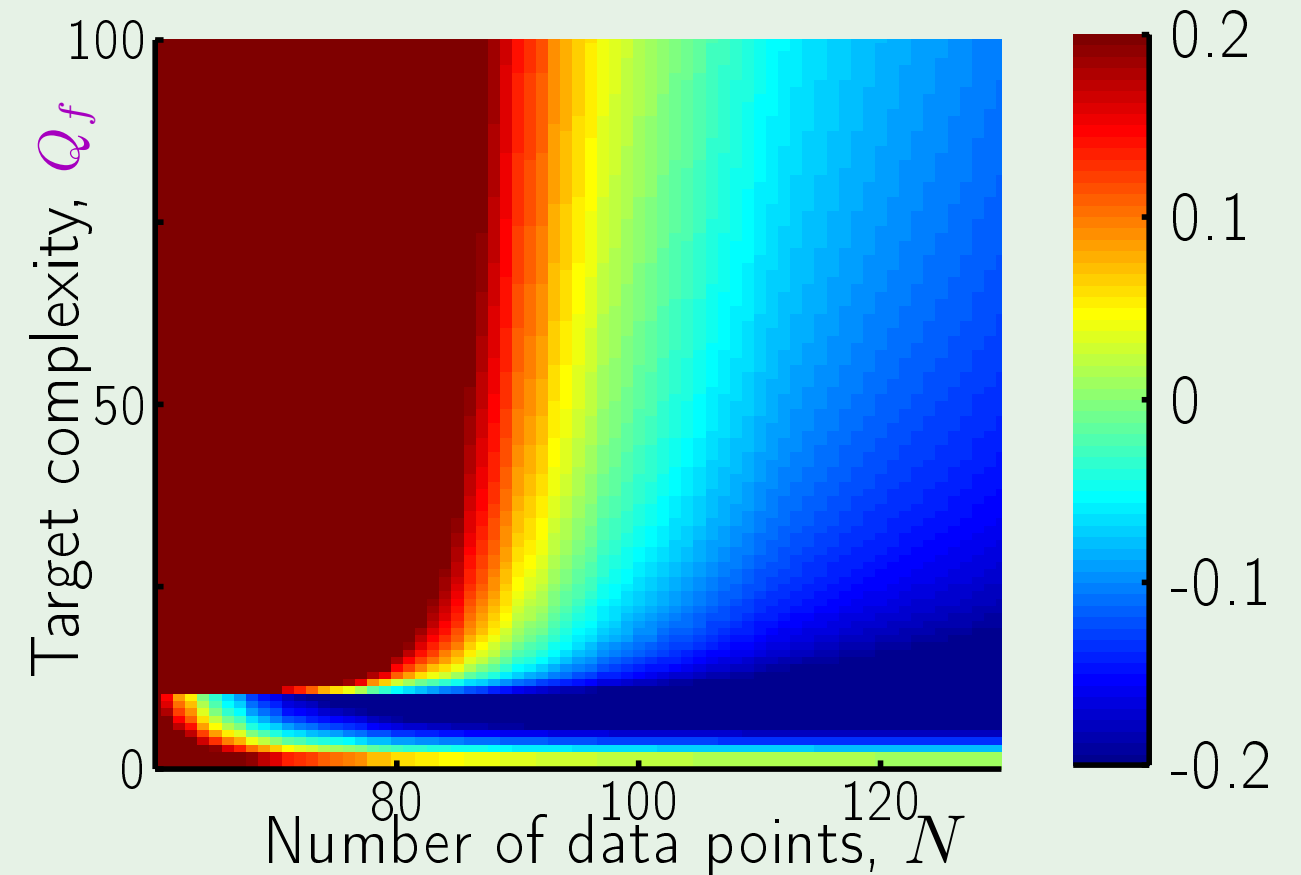


# Impact on overfitting

Deterministic noise and  $Q_f$

Finite  $N$ :  $\mathcal{H}$  tries to fit the noise

The finite dataset means that the hypothesis set may be able to capture a small amount of stochastic or deterministic noise. With more data points, the model is less affected by stochastic noise as it cannot fit it. With a smaller dataset, our model gains the unfortunate ability of being able to fit both the stochastic and deterministic noise, and so the model will fit it - even if it is deterministic where there is no point in fitting it because we know that our model is out of our ability and we have no way to generalize out of sample for it



how much overfit

# Noise and bias-variance

Recall the decomposition:

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

What if  $f$  is a noisy target?

$$y = f(\mathbf{x}) + \epsilon(\mathbf{x}) \quad \mathbb{E} \left[ \epsilon(\mathbf{x}) \right] = 0$$

## A noise term

$$\begin{aligned}\mathbb{E}_{\mathcal{D}, \epsilon} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - y \right)^2 \right] &= \mathbb{E}_{\mathcal{D}, \epsilon} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}) - \epsilon(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}, \epsilon} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) + \bar{g}(\mathbf{x}) - f(\mathbf{x}) - \epsilon(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E}_{\mathcal{D}, \epsilon} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 + \left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 + \left( \epsilon(\mathbf{x}) \right)^2 \right. \\ &\quad \left. + \text{cross terms} \right]\end{aligned}$$

We then take the expectation value w.r.t.  $\mathbf{x}$  (as in the previous decomposition) to get the variance and bias independent of  $\mathbf{x}$

## Actually, two noise terms

Last term characterizes the distance from  $f$ , the target proper, to the actual output (which has a noise aspect to it)

$$\dots \underbrace{\mathbb{E}_{\mathcal{D}, \mathbf{x}} \left[ \left( g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) \right)^2 \right]}_{\text{var}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[ \left( \bar{g}(\mathbf{x}) - f(\mathbf{x}) \right)^2 \right]}_{\substack{\text{bias} \\ \uparrow \\ \text{deterministic noise}}} + \underbrace{\mathbb{E}_{\epsilon, \mathbf{x}} \left[ \left( \epsilon(\mathbf{x}) \right)^2 \right]}_{\substack{\sigma^2 \\ \uparrow \\ \text{stochastic noise}}}$$

$\sigma^2$  is the "energy" of the noise

The bias characterizes the ability of our best possible hypothesis to approximate  $f$ , so this is the energy of the deterministic noise.

This presentation of the variance and two types of noise allows us to treat the noises in the same way - increasing the number of examples gives us a better variance (the red region representing the variance shrinks), the stochastic and deterministic noise are both inevitable (there is nothing we can do about stochastic noise, and the deterministic noise is fixed given a hypothesis set). Note that for a given dataset, both noises have a finite version on the data points and the algorithm will try to fit them, so both noises will affect the variance by making the fit more susceptible to choosing different hypotheses depending on the dataset (since each noise will affect each dataset differently) - not because it is indicated by the target function we want to learn but because there is a noise present in the sample, which we blindly follow because we can't distinguish noise from signal, and so we end up with more variety/possible choices of  $g$ , so larger red region, so worse variance and overfit.

# Outline

- What is overfitting?
- The role of noise
- Deterministic noise
- Dealing with overfitting

# Two cures

**Regularization:**

Putting the brakes

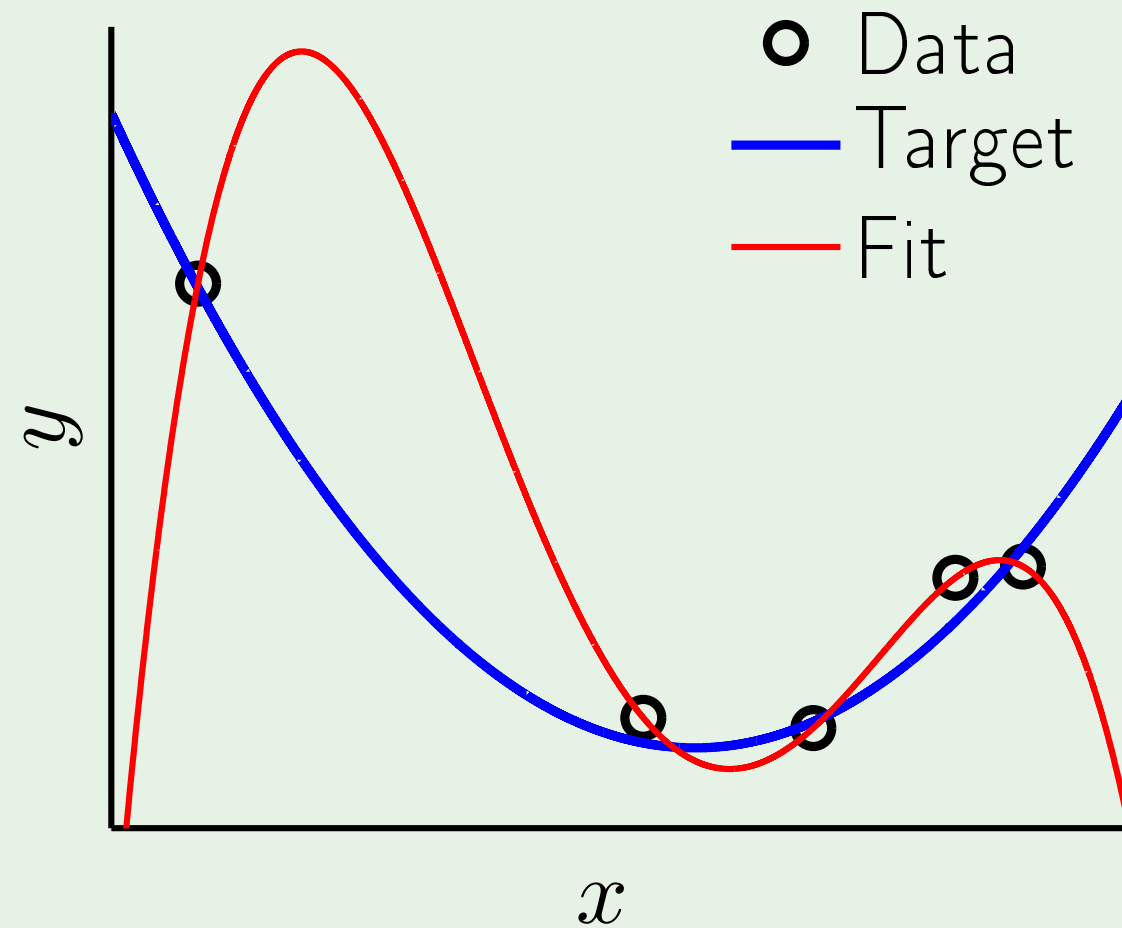
stop further fitting/minimization of  $E_{in}$  to prevent  $E_{out}$  becoming worse

**Validation:**

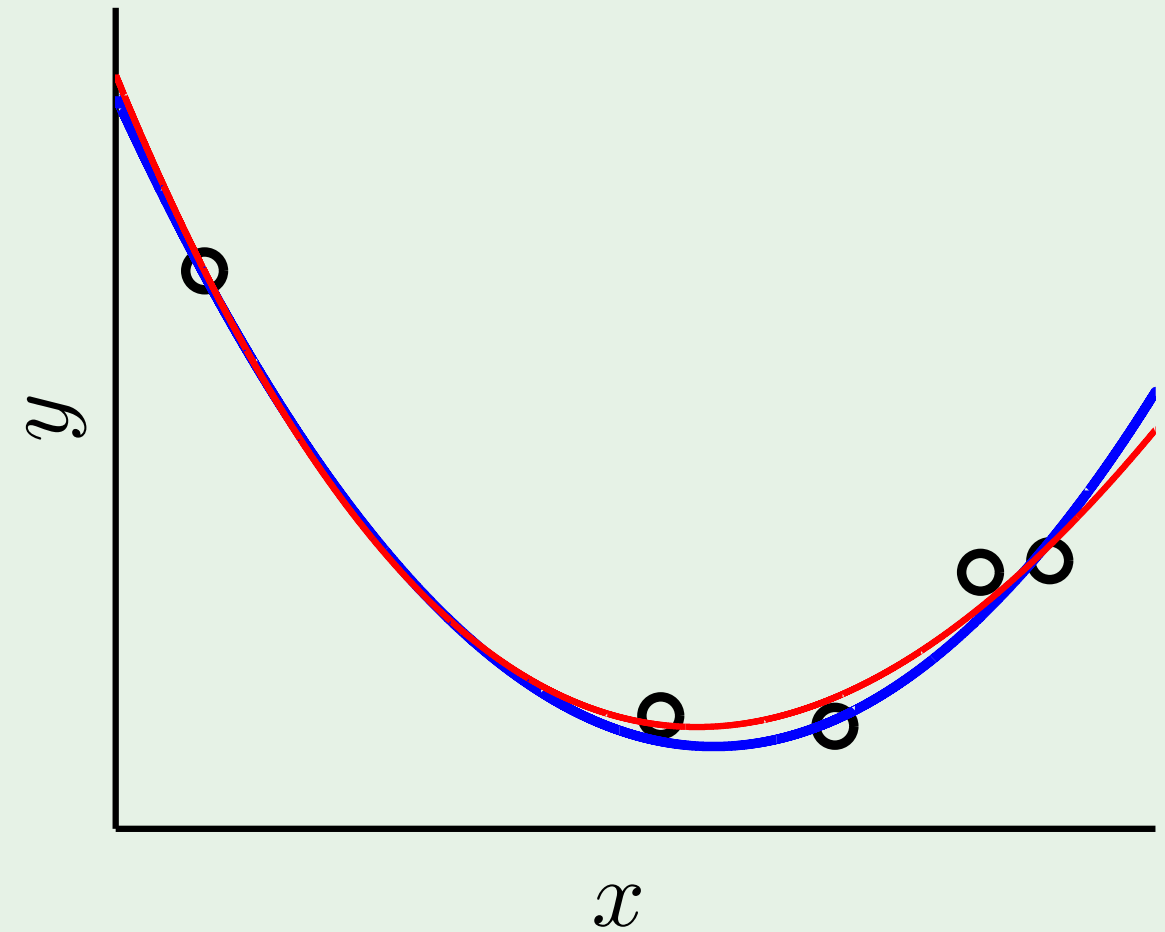
Checking the bottom line

Since  $E_{in}$  is not a very good indicator of generalization, we want another way of checking the generalization ability/what is happening out of sample - so we can avoid overfitting by checking what is happening in the quantity we care about ( $E_{out}$ )

# Putting the brakes



free fit



restrained fit