

Bias-Variance Tradeoff

Overview of the second half the semester

1. A little bit Learning Theory

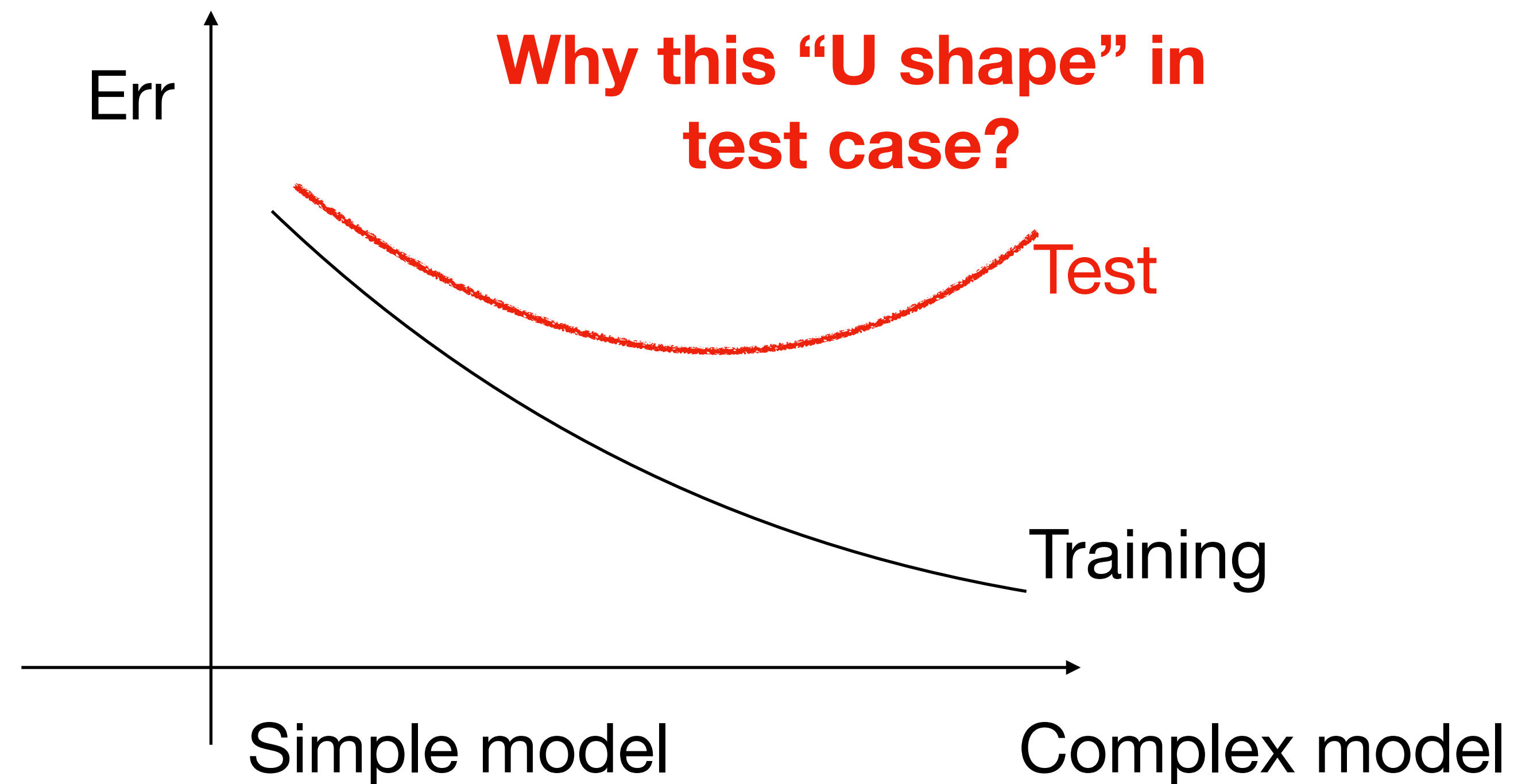
2. Make our linear models nonlinear (Kernel)

3. How to combine multiple classifiers into a stronger one (Bagging & Boosting)?

4. Intro of Neural Networks (old and new)

Objective

Understand Bias-Variance tradeoff — When and why your ML models work (or don't work)

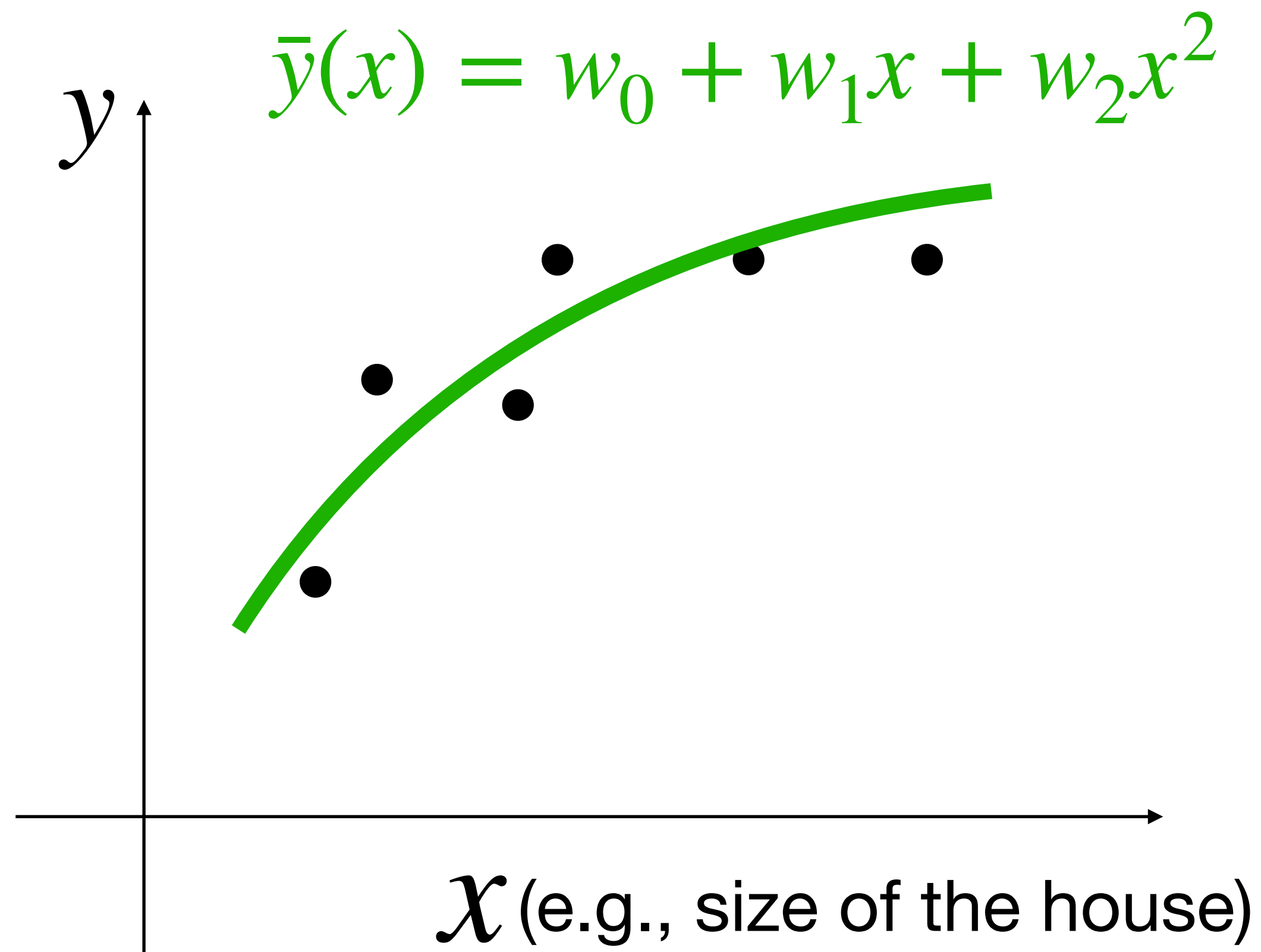


Outline of Today

1. Intro on Underfitting/Overfitting and Bias/Variance
2. Derivation of the Bias-Variance Decomposition

Bayes optimal predictor

Consider regression problem w/ dataset $\mathcal{D} = \{x, y\}, (x, y) \sim P, x \in \mathbb{R}, y \in \mathbb{R}$

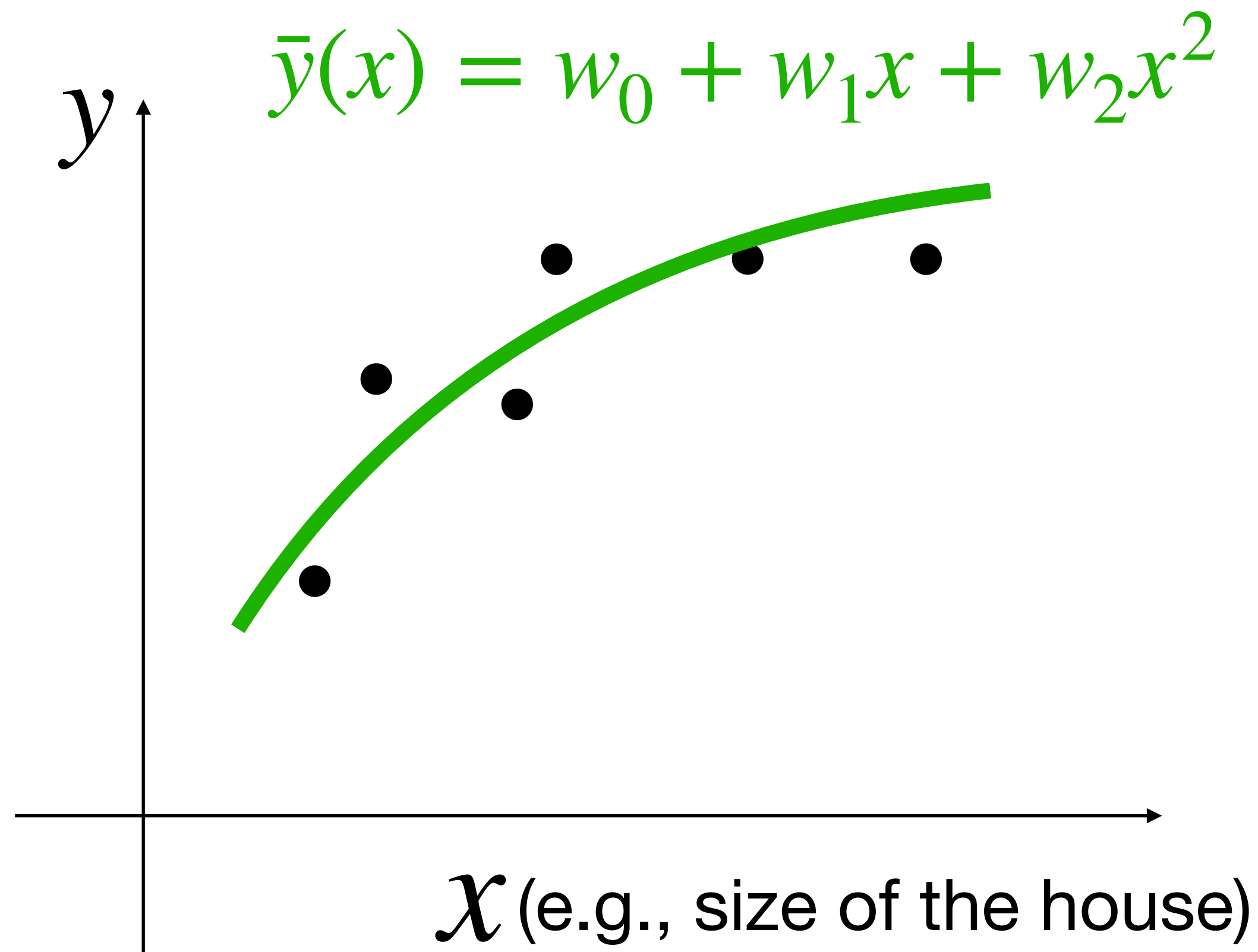


The Bayes optimal regressor:

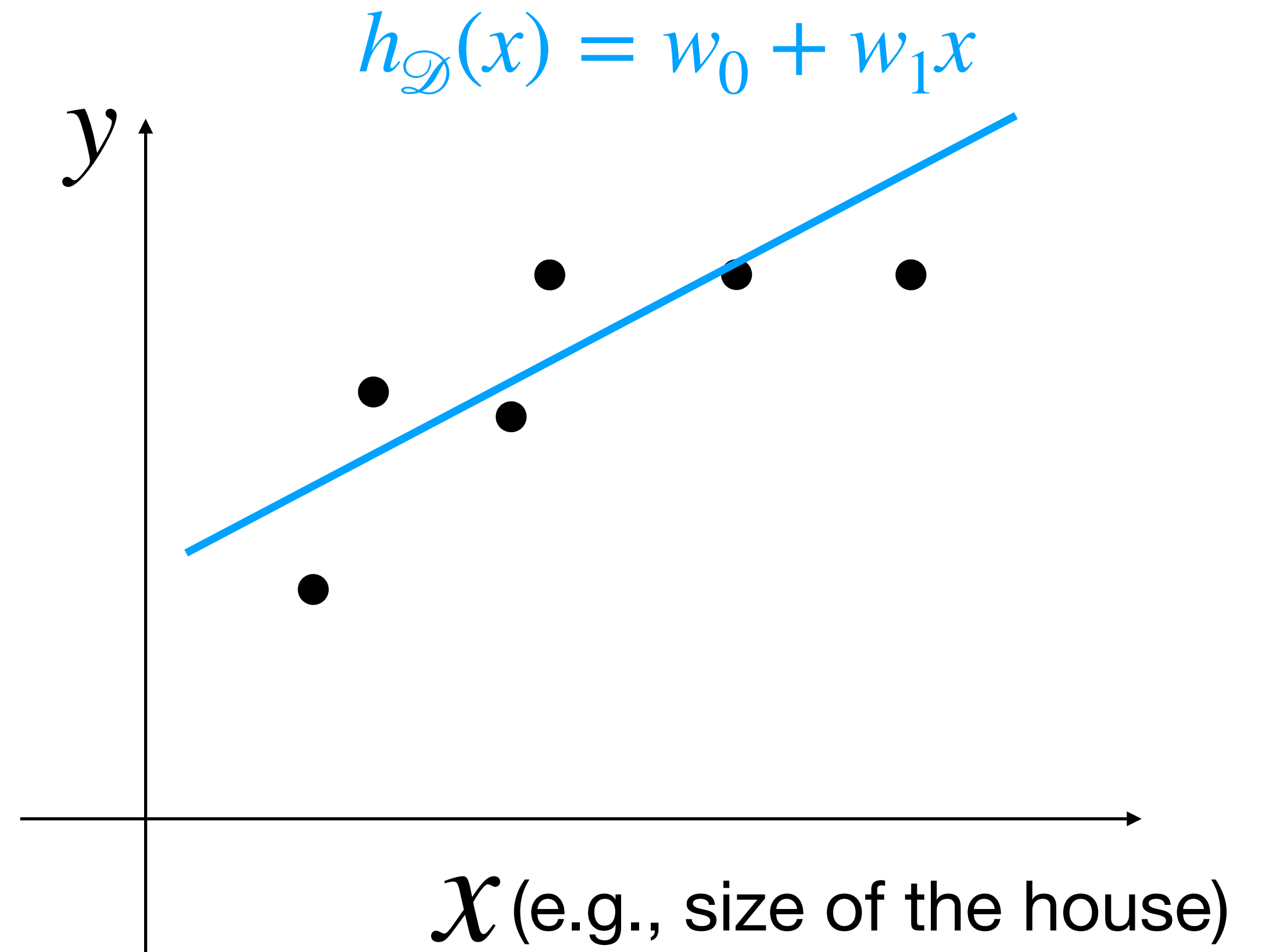
$$\bar{y}(x) := \mathbb{E}[y \mid x]$$

The best we could do, cannot beat this one

Underfitting



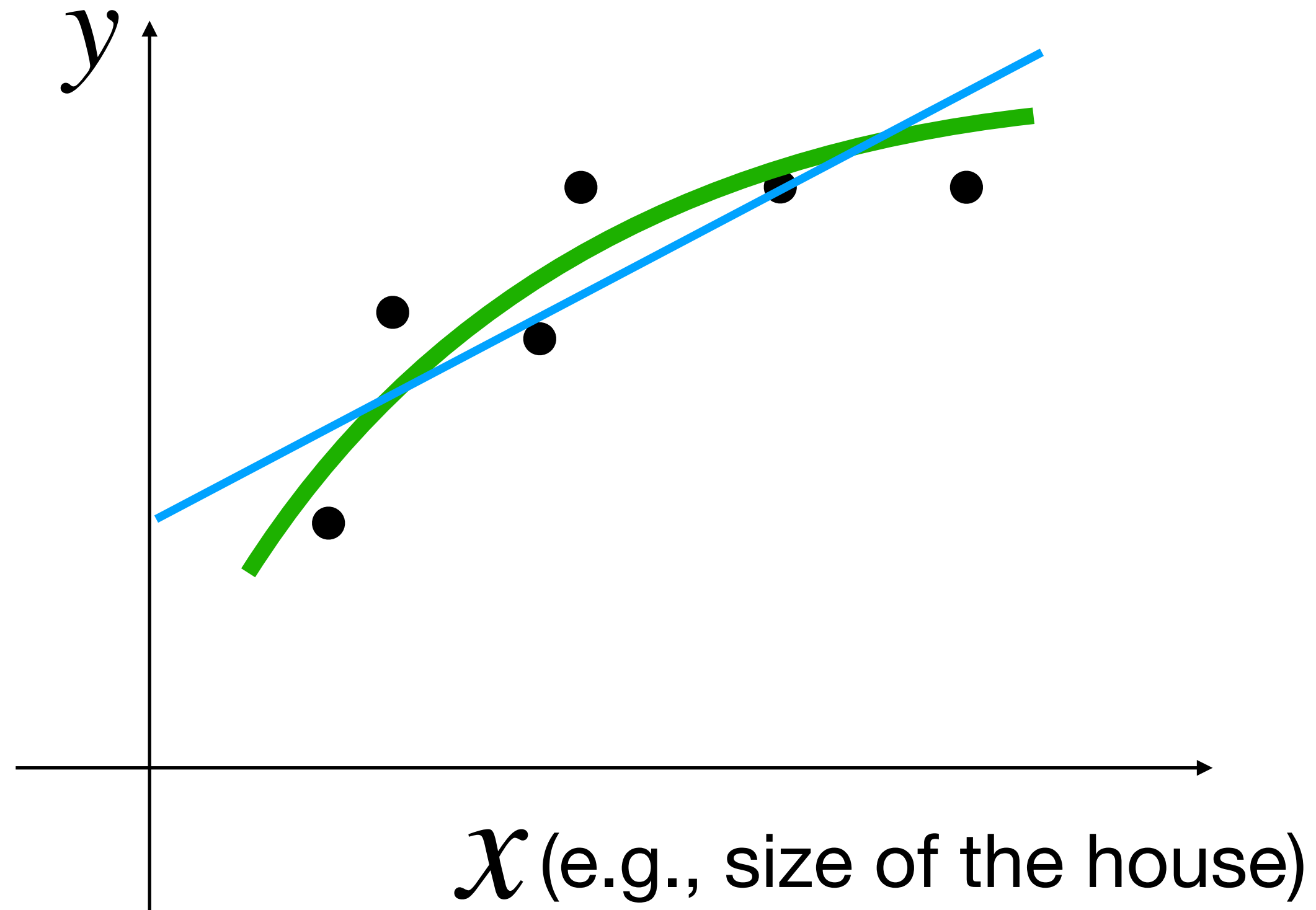
(Just right)



Underfitting

Underfitting

Just right versus Underfitting

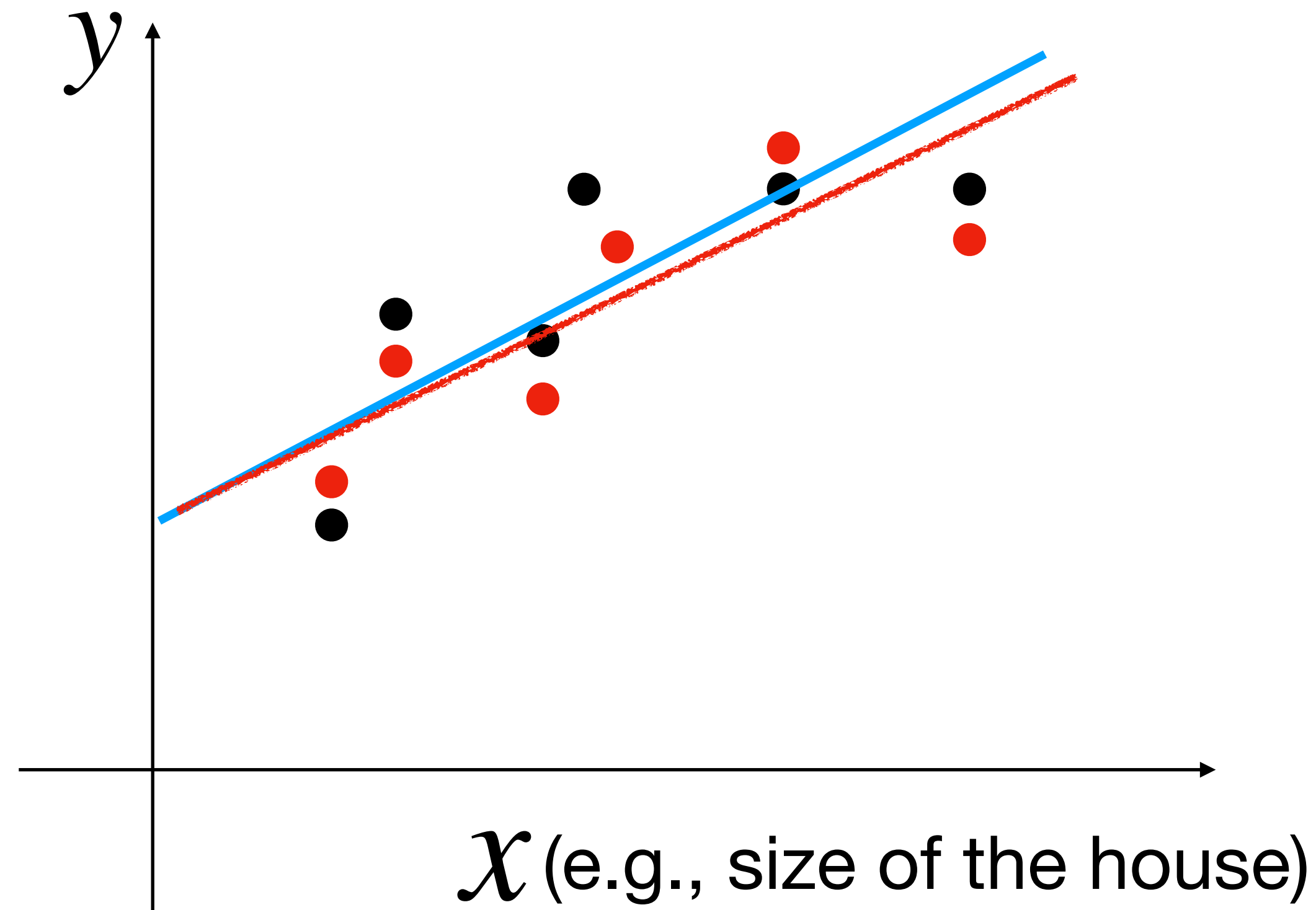


Bias:

Bias towards to linear models

Underfitting

Now let's redo linear regression on a **different dataset \mathcal{D}'** (but from the same distribution)



The new linear function does not differ too much from the old one

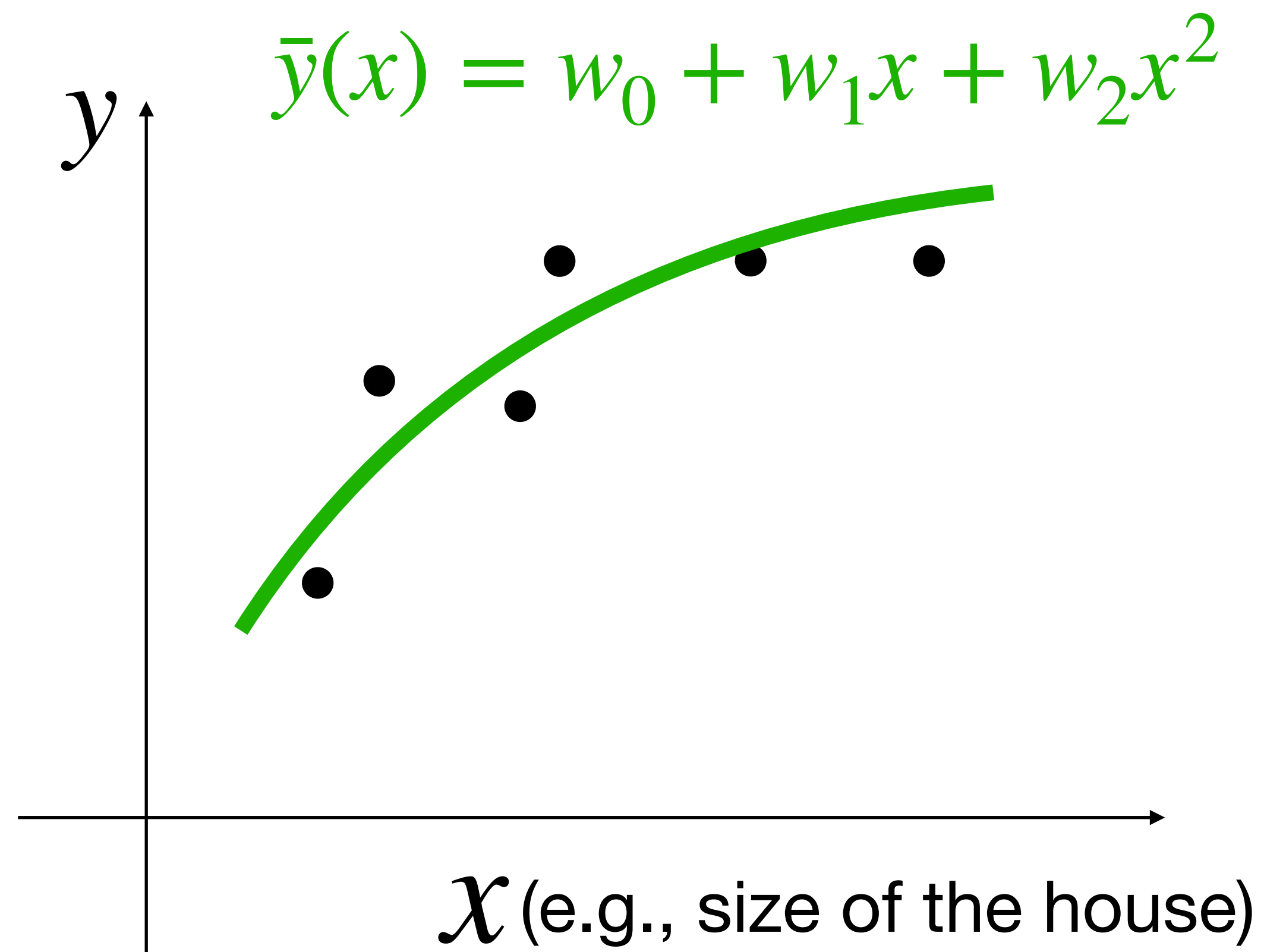
This is called low variance

Q: what happens when our linear predictor is $h(x) = w_0$?

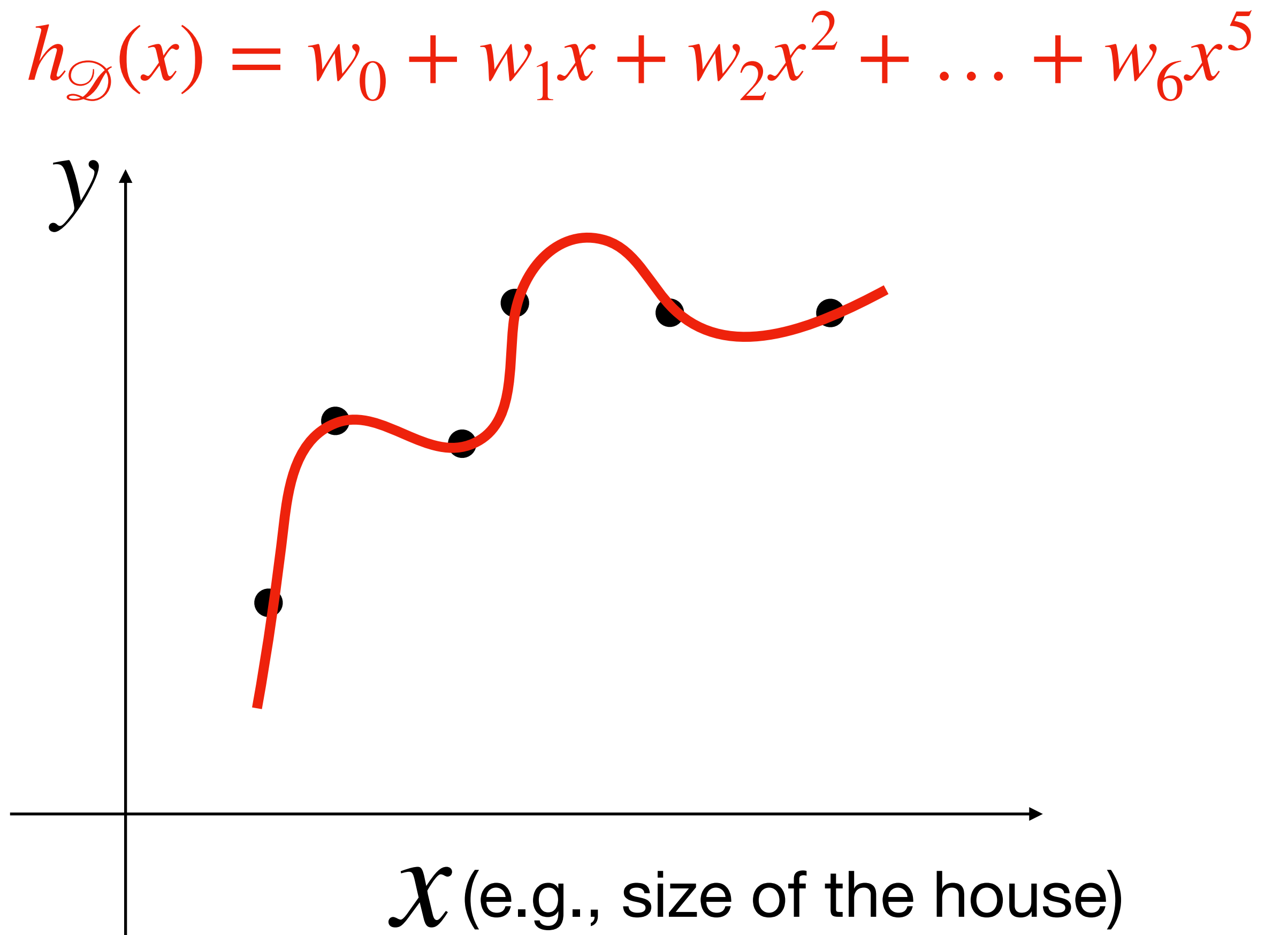
Summary on underfitting

1. Often our model is too simple, i.e., we bias towards too simple models
2. This causes underfitting, i.e., we cannot capture the trend in the data
3. In this case, we have large bias, but low variance (think about the $h(x) = w_0$ case)

Overfitting



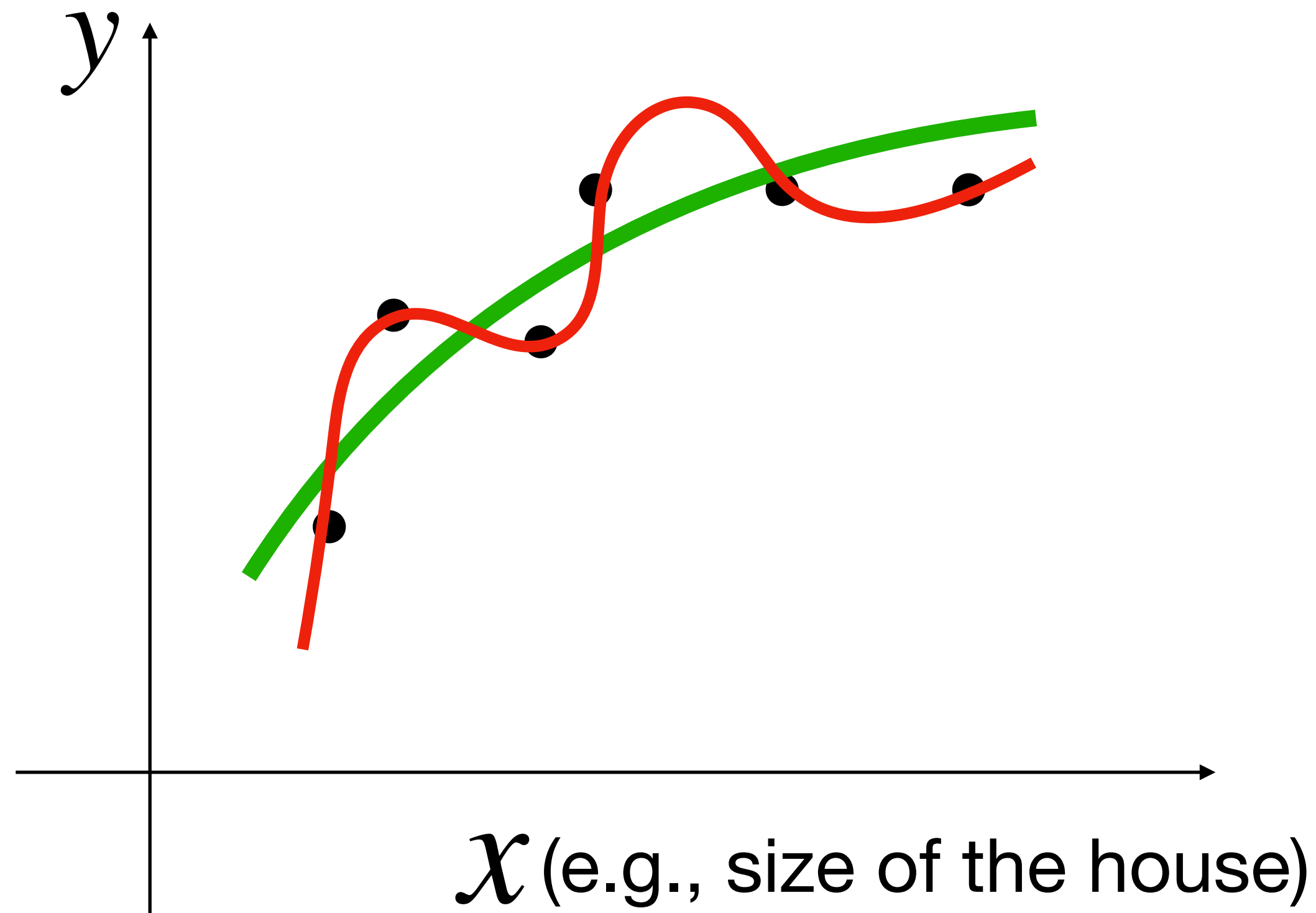
(Just right)



Overfitting

Overfitting

Just right versus Overfitting



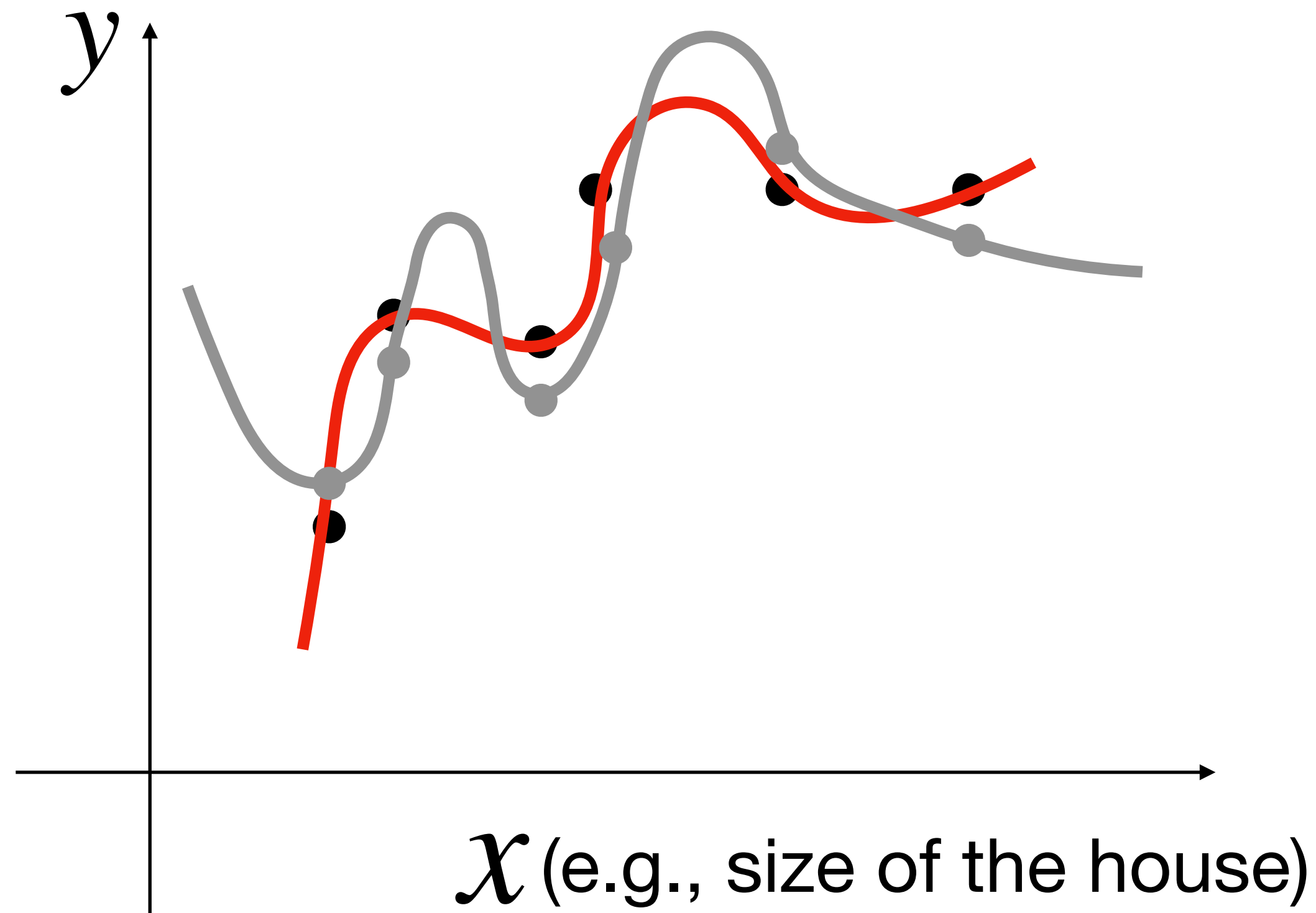
No strong bias:

Our hypothesis class is all
polynomials up to 5-th order

i.e., a priori, no strong bias towards
linear or quadratic, or cubic, etc

Overfitting

Redo the higher-order polynomial fitting on different dataset \mathcal{D}'



The new function can differ a lot from the old one

This is called high variance

Summary on Overfitting

1. Often our model is too complex (e.g., can fit noise perfectly to achieve zero training error)
 2. This causes overfitting, i.e., cannot generalize well on unseen test example
 3. In this case, we have small bias, but large variance
(tiny change on the dataset cause large change in the fitted functions)

Outline of Today

1. Intro on Underfitting/Overfitting and Bias/Variance

2. Derivation of the Bias-Variance Decomposition

Generalization error

Given dataset \mathcal{D} , a hypothesis class \mathcal{H} , squared loss $\ell(h, x, y) = (h(x) - y)^2$,
denote $h_{\mathcal{D}}$ as the ERM solution

We are interested in the generalization error of $h_{\mathcal{D}}$:

$$\mathbb{E}_{\mathcal{D}} \mathbb{E}_{x, y \sim P} (h_{\mathcal{D}}(x) - y)^2$$

Q: how to estimate this in practice?

The expectation of our model $h_{\mathcal{D}}$

Since $h_{\mathcal{D}}$ is random, we consider its expected behavior:

$$\bar{h} := \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}]$$

In other words, we have:

$$\bar{h}(x) = \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}(x)], \forall x$$

Q: what is \bar{h} is the case where hypothesis is $h(x) = w_0$?

$$\text{A: } \bar{h}(x) = \mathbb{E}_y[y]$$

Formal definition of Bias and Variance

$$\bar{h} := \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}] \quad \bar{y}(x) := \mathbb{E}[y | x]$$

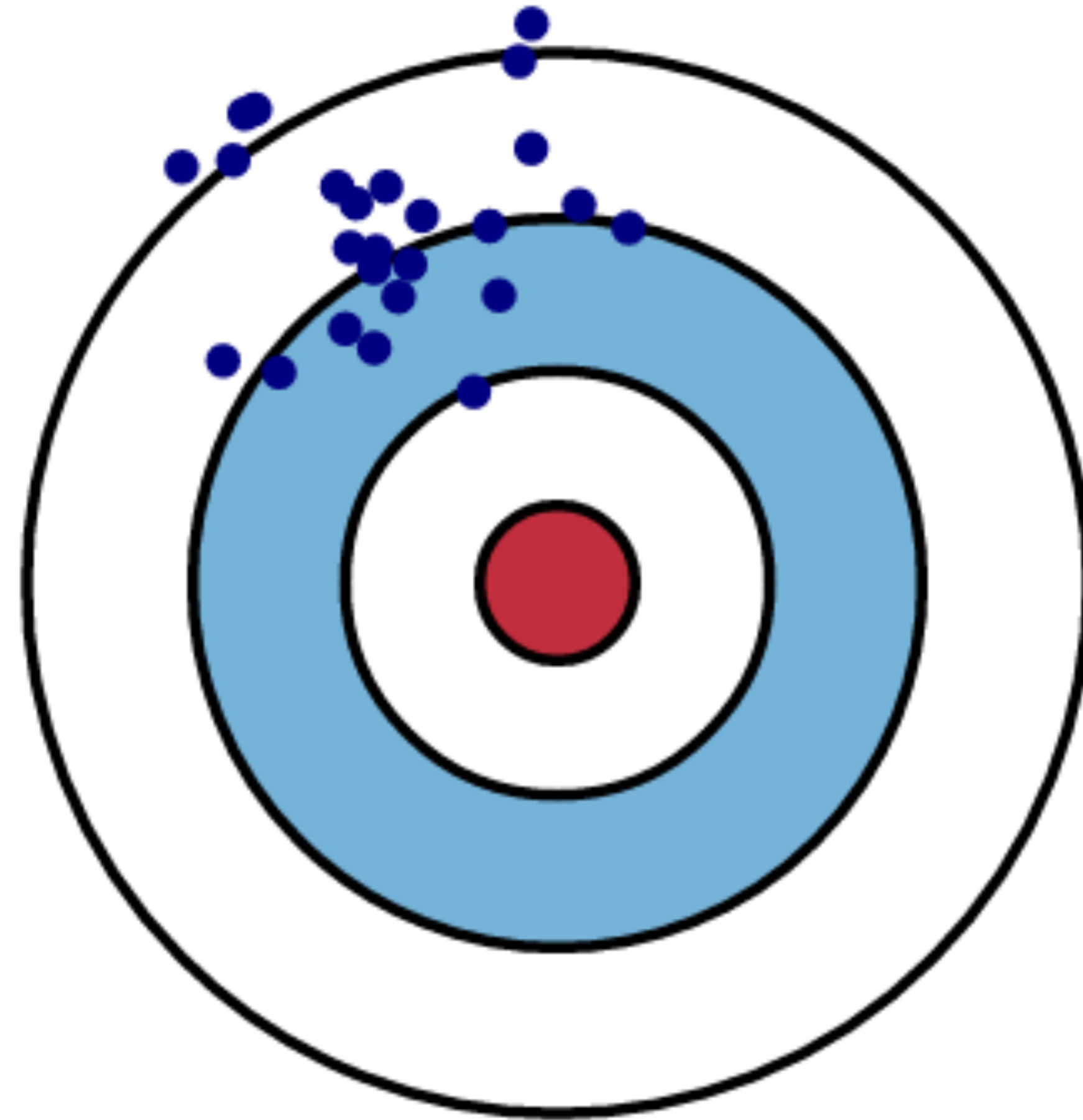
Bias²: (squared) difference between \bar{h} and the best $\bar{y}(x)$, i.e., $\mathbb{E}_x (\bar{y}(x) - \bar{h}(x))^2$

Difference between our mean and the best

Variance: difference from \bar{h} and $h_{\mathcal{D}}$, i.e., $\mathbb{E}_{\mathcal{D}} \mathbb{E}_x (h_{\mathcal{D}}(x) - \bar{h}(x))^2$

Fluctuation of our random model around its mean

Bias-Variance illustration



Generalization error decomposition

$$\bar{h} := \mathbb{E}_{\mathcal{D}} [h_{\mathcal{D}}] \quad \bar{y}(x) := \mathbb{E}[y \mid x]$$

What we gonna show now:

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} \mathbb{E}_{x,y \sim P} (h_{\mathcal{D}}(x) - y)^2 \\ &= \mathbf{Bias}^2 + \mathbf{Variance} + \text{Noise (unavoidable, independent of Algs/models)} \end{aligned}$$

We will use the following trick twice: $(x - y)^2 = (x - z)^2 + (z - y)^2 + 2(x - z)(z - y)$

$$\begin{aligned}
& \mathbb{E}(h_{\mathcal{D}}(x) - y)^2 \\
&= \mathbb{E}(h_{\mathcal{D}}(x) - \bar{h}(x) + \bar{h}(x) - y)^2 \\
&= \mathbb{E}(h_{\mathcal{D}}(x) - \bar{h}(x))^2 + \mathbb{E}(\bar{h}(x) - y)^2 - 2\mathbb{E}_{\mathcal{D},x,y} [(h_{\mathcal{D}}(x) - \bar{h}(x))(\bar{h}(x) - y)]
\end{aligned}$$

This term is zero since:

$$\begin{aligned}
& \mathbb{E}_{x,y,\mathcal{D}} [(h_{\mathcal{D}}(x) - \bar{h}(x))(\bar{h}(x) - y)] \\
&= \mathbb{E}_{x,y} [\mathbb{E}_{\mathcal{D}}(h_{\mathcal{D}}(x) - \bar{h}(x)) \cdot (\bar{h}(x) - y)] \\
&= \mathbb{E}_{x,y} [(\bar{h}(x) - \bar{h}(x)) \cdot (\bar{h}(x) - y)]
\end{aligned}$$

$$\begin{aligned}
 & \mathbb{E}(h_{\mathcal{D}}(x) - y)^2 \\
 = & \underbrace{\mathbb{E}(h_{\mathcal{D}}(x) - \bar{h}(x))^2}_{\text{Variance}} + \mathbb{E}(\bar{h}(x) - y)^2 \\
 = & \mathbb{E}(\bar{h}(x) - \bar{y}(x) + \bar{y}(x) - y)^2 \\
 = & \mathbb{E}(\bar{h}(x) - \bar{y}(x))^2 + \mathbb{E}(\bar{y}(x) - y)^2 \\
 & + 2\mathbb{E}(\bar{h}(x) - \bar{y}(x))(\bar{y}(x) - y)
 \end{aligned}$$

This term is zero since:

$$\begin{aligned}
 & = \mathbb{E}_x \left[(\bar{h}(x) - \bar{y}(x)) \cdot \mathbb{E}_{y|x}(\bar{y}(x) - y) \right] \\
 & = \mathbb{E}_x \left[(\bar{h}(x) - \bar{y}(x)) \cdot (\bar{y}(x) - \mathbb{E}_{y|x}[y]) \right]
 \end{aligned}$$

Putting the derivations together, we arrive at:

$$\mathbb{E}(h_{\mathcal{D}}(x) - y)^2 = \underbrace{\mathbb{E}(h_{\mathcal{D}}(x) - \bar{h}(x))^2}_{\text{Variance}} + \underbrace{\mathbb{E}(\bar{h}(x) - \bar{y}(x))^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}(\bar{y}(x) - y)^2}_{\text{Noise}}$$

Note that the noise term is independent of training algorithms / models