

# Bias-Variance Decomposition

Machine Learning- CS-433

1 Oct 2025

Robert West

(Slide credits: Martin Jaggi & Nicolas Flammarion)



# Last time

How can we judge if a given predictor is good?

How to select the best models of a family?

- Bound the difference between the true and empirical risks
- Split data into train and test sets (learn with the train and test on the test)

Motivation: Hyperparameter search (where hyperparameters often control model complexity)

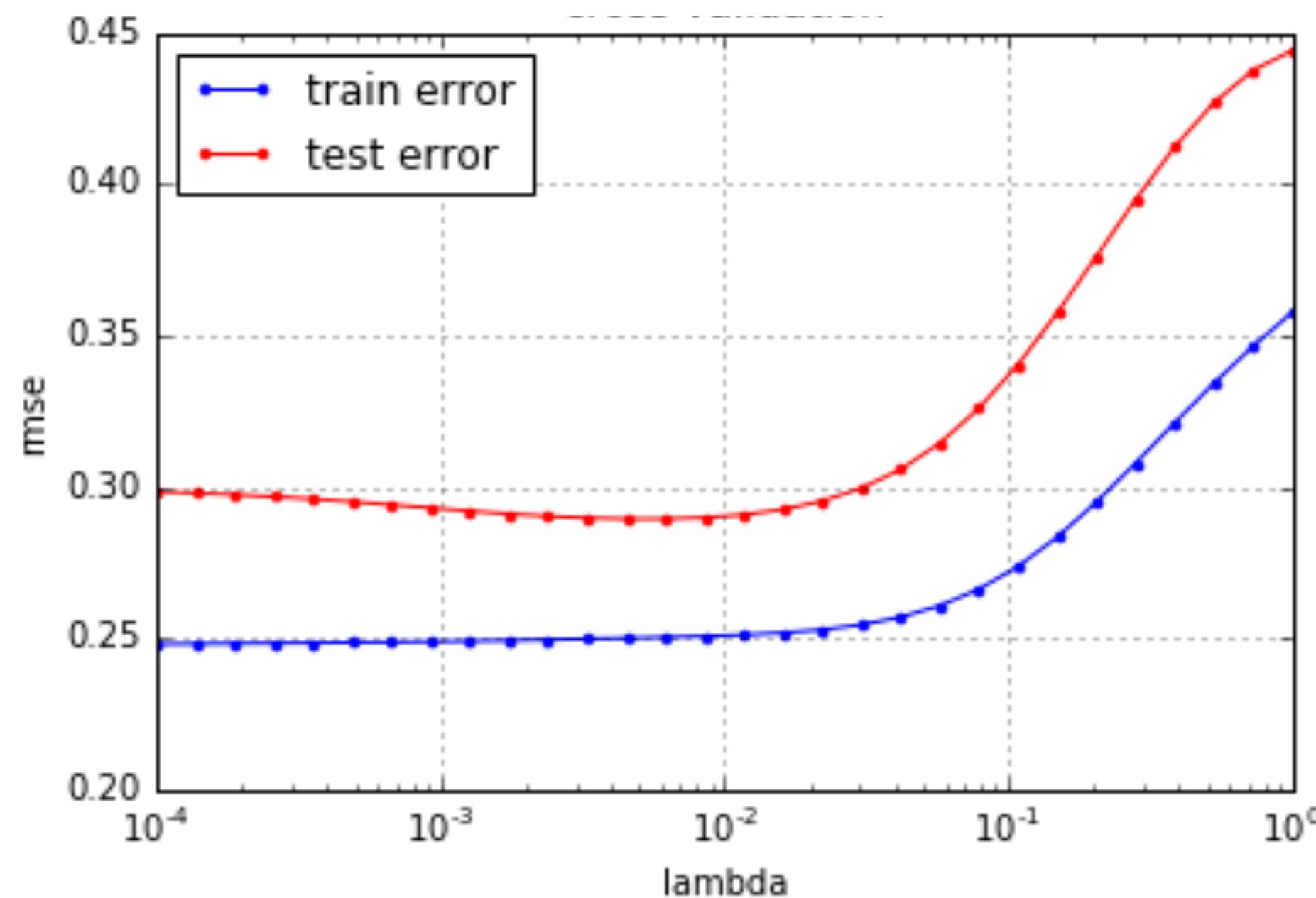
But we haven't investigated the role of the complexity of the class

C

“

“

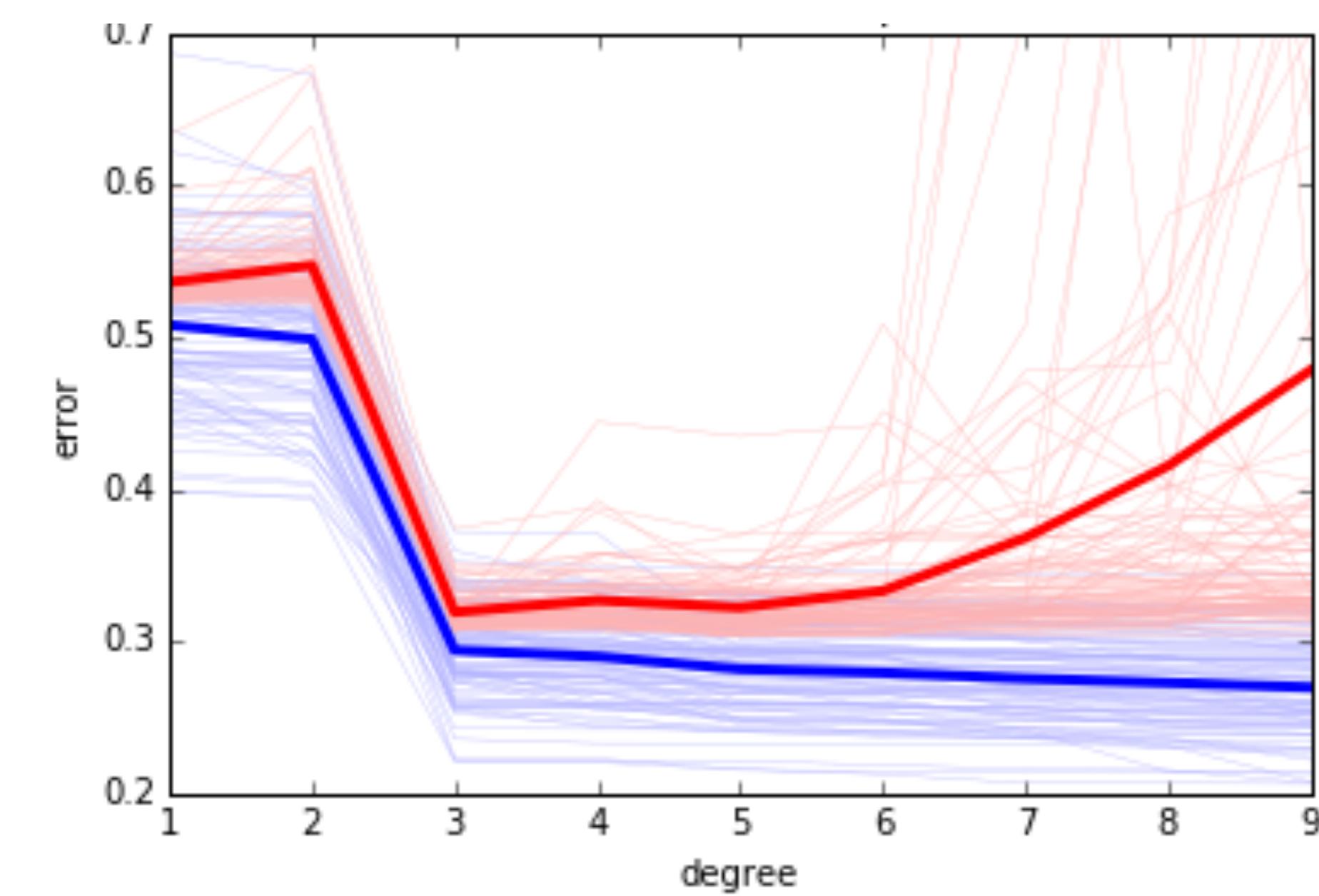
# Model selection curves



Ridge regression

$\lambda$

less complexity

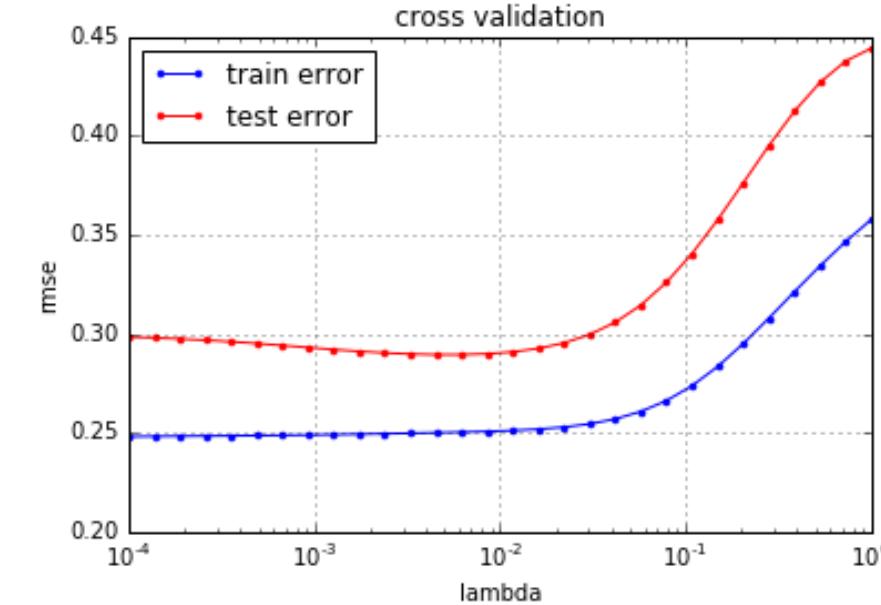


Degree in case of a polynomial feature expansion

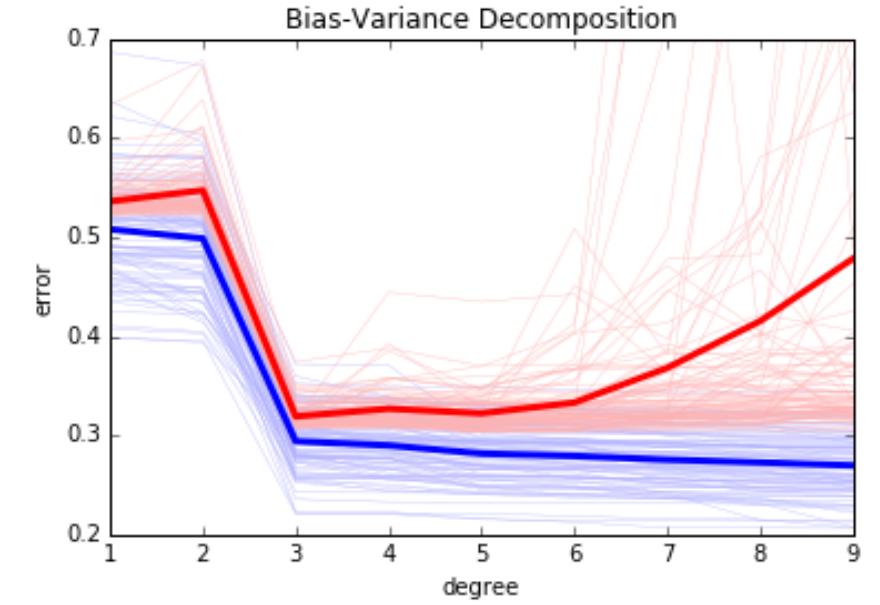
$\leftarrow$

less complexity

# Today



Ridge regression

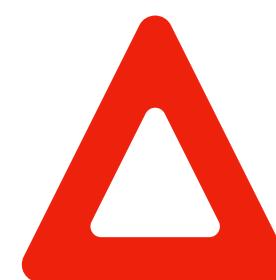


Polynomial feature  
expansion

How does the risk behave as a function of the complexity of the model class?

→ ***Bias-Variance tradeoff***

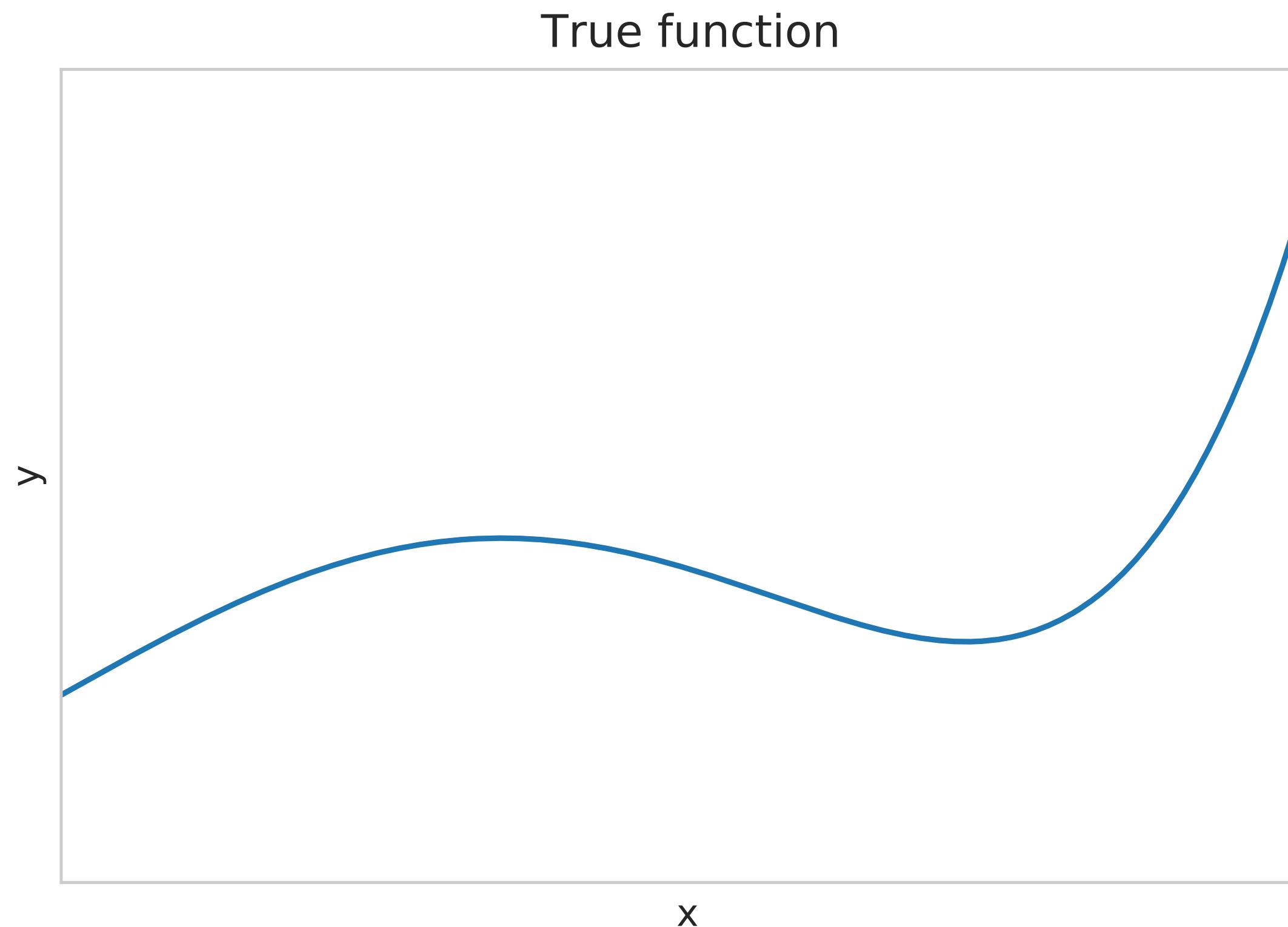
It will help us to decide how complex and rich we should make our model



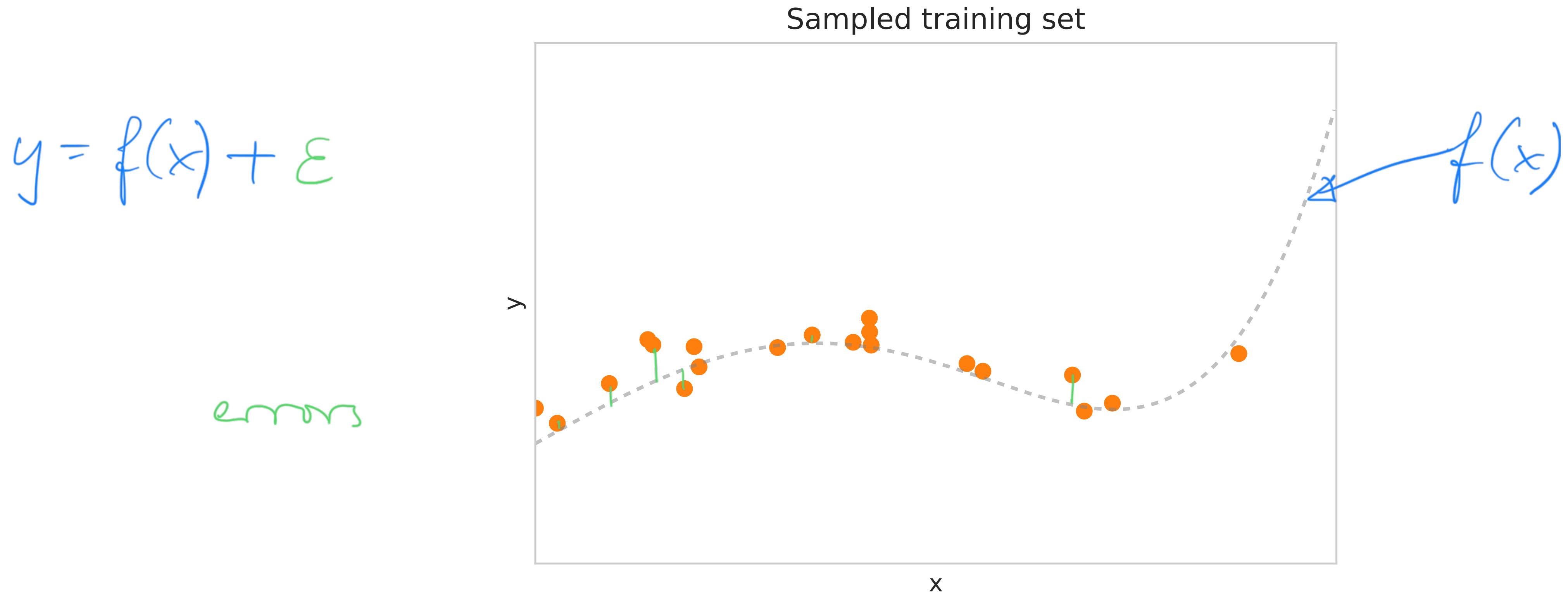
Before: quantitative

Now: ***qualitative***

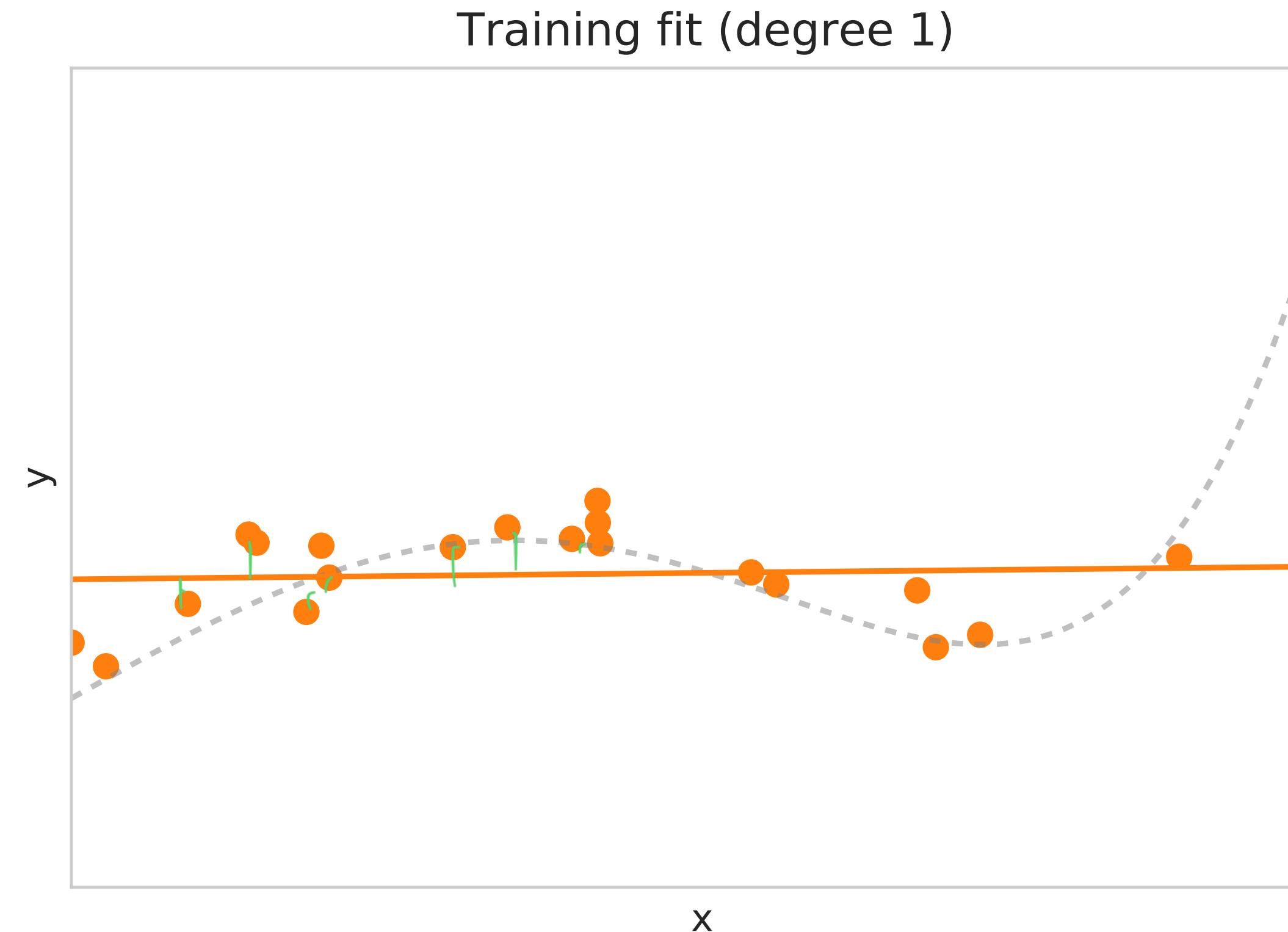
# A small experiment: 1D-regression



# A small experiment: 1D-regression

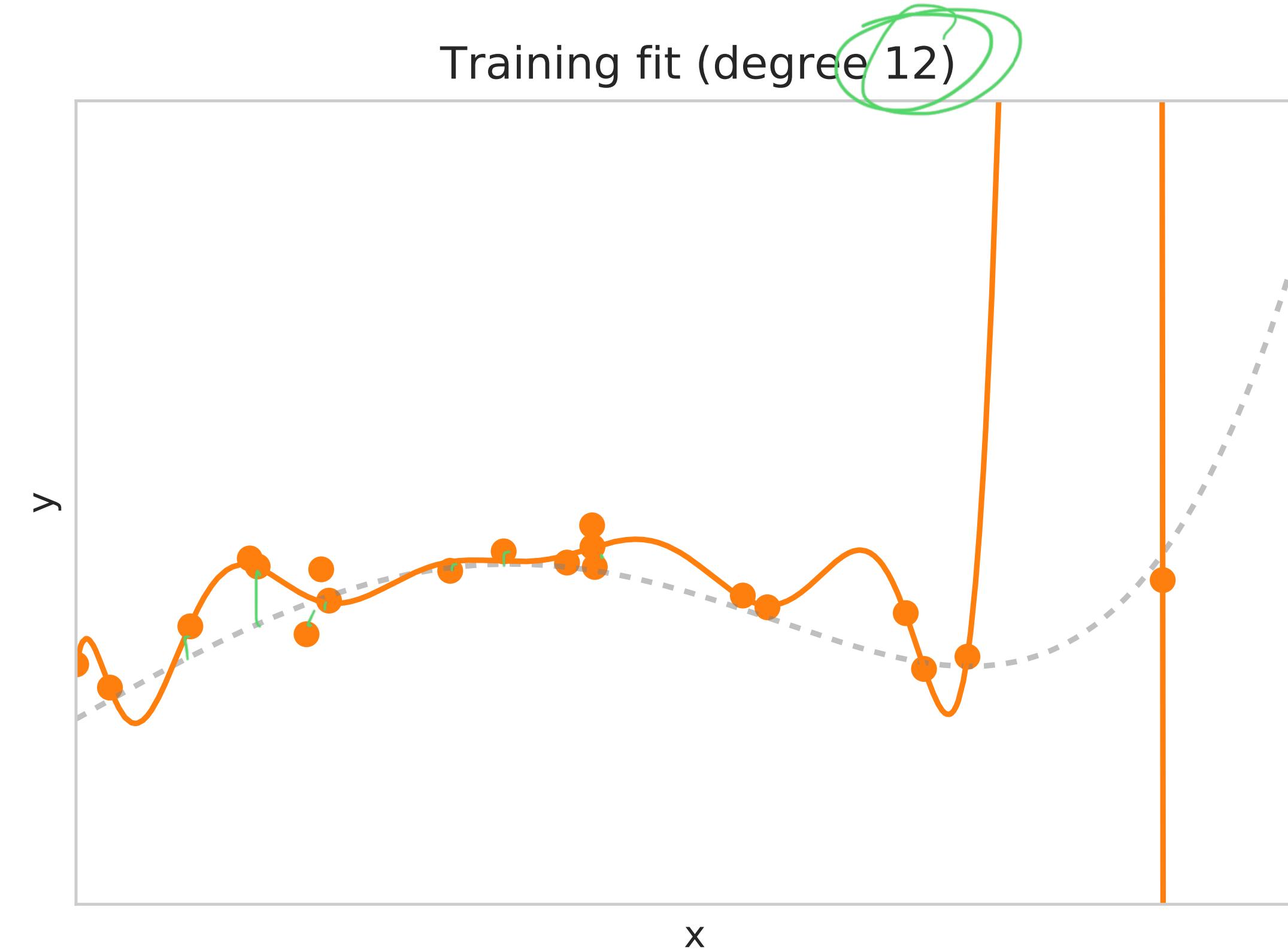


# Simple model: bad fit



No linear function would be a good predictor. The model class is not rich enough

# Complex model: good fit?

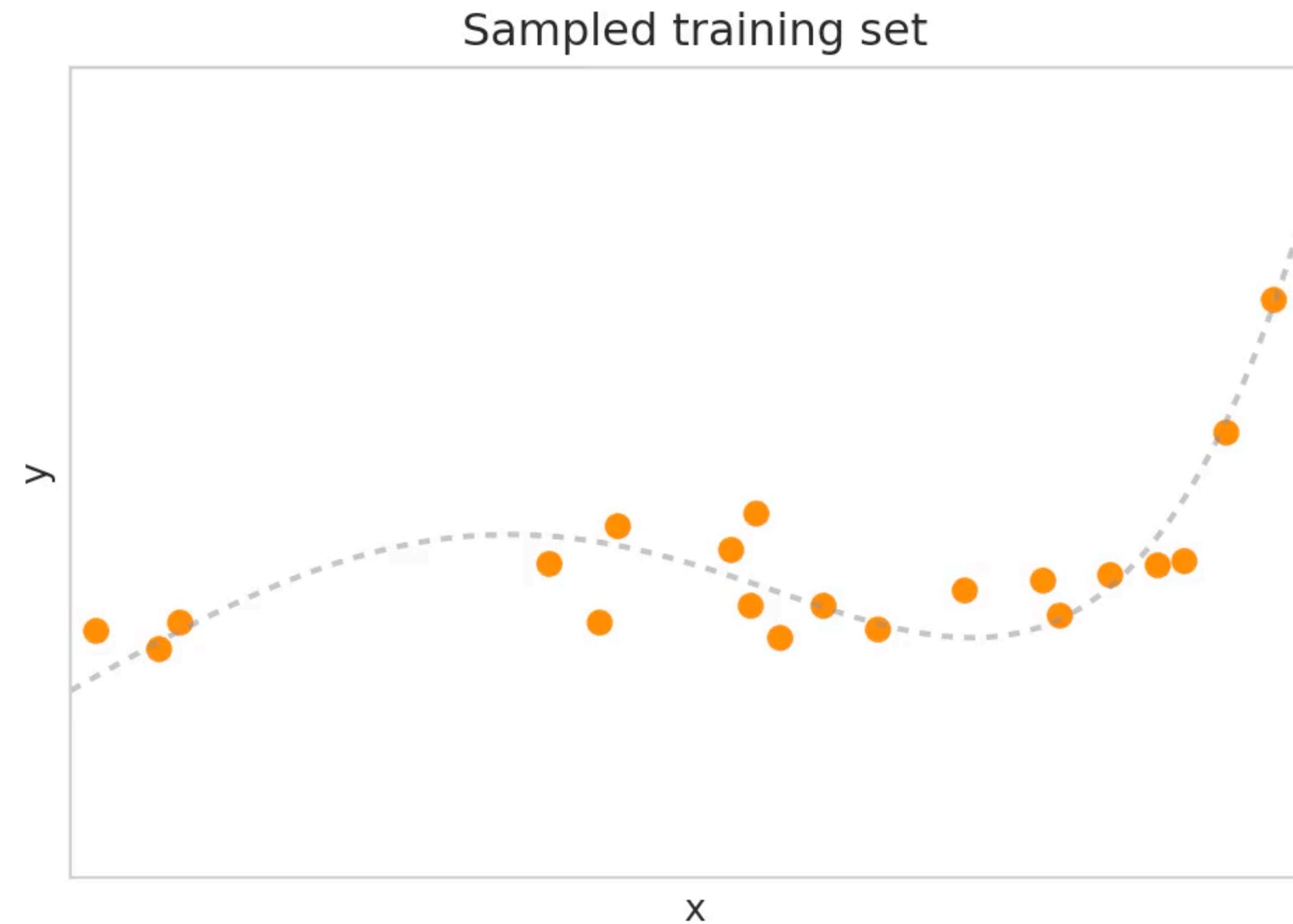


Linear regression using polynomial feature expansion  $(x, x^2, x^3, \dots, x^d)$

The maximum degree  $d$  measures the complexity of the class

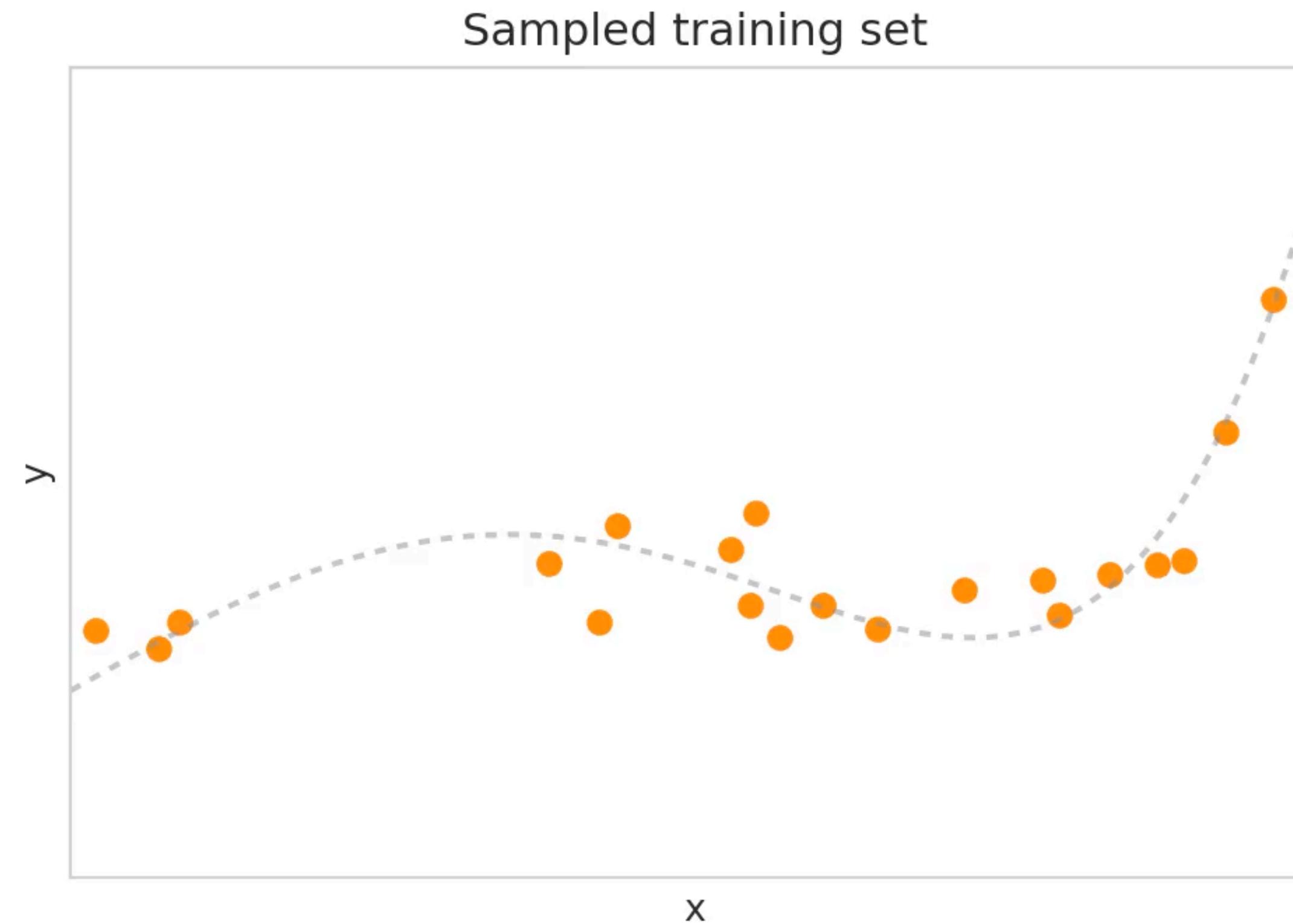
High-degree polynomial (complex model) will be a good fit. But?

# But there is randomness in the data



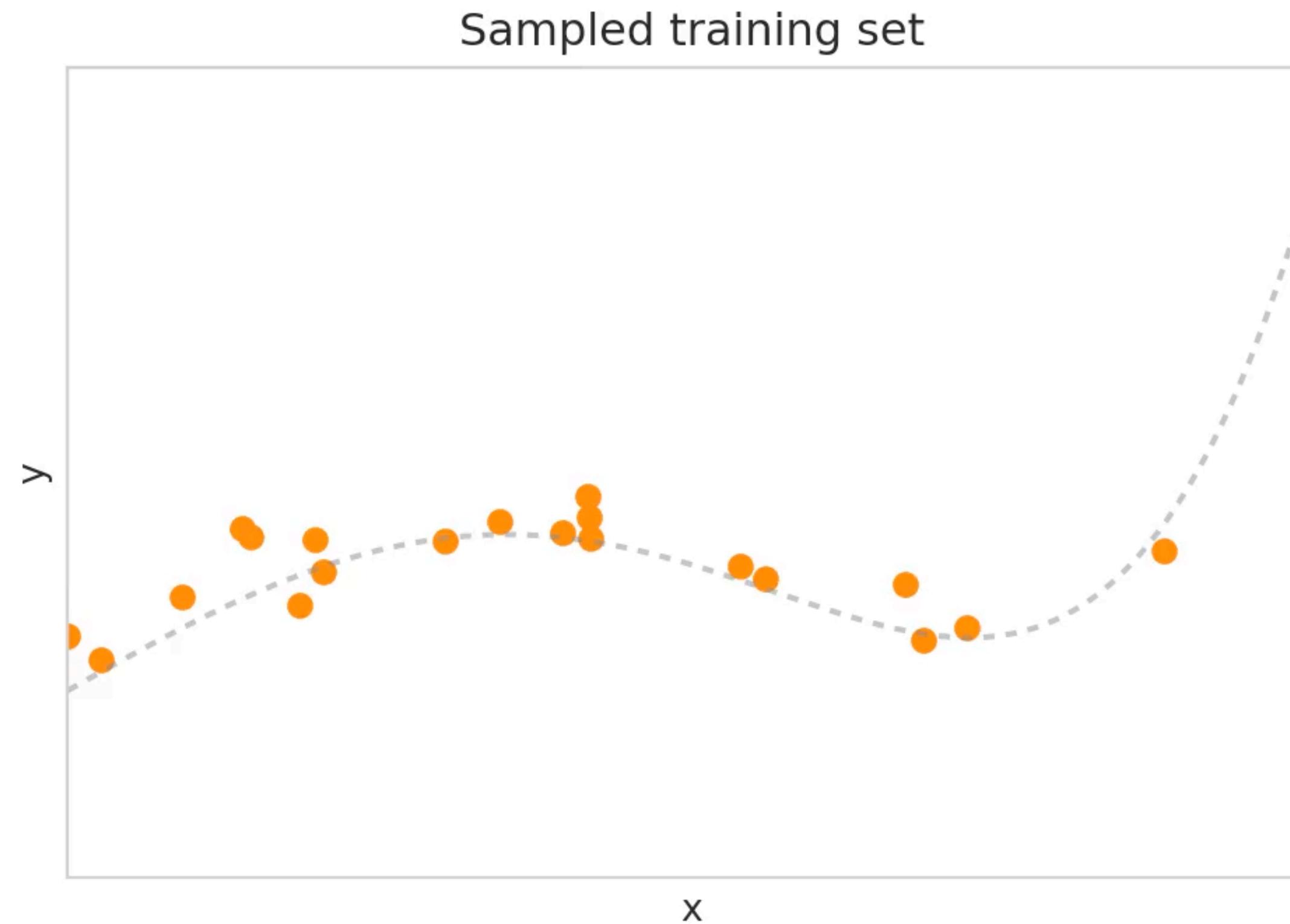
We have observed one particular  $S_{\text{train}}$  but we could have observed several others!

# But there is randomness in the data



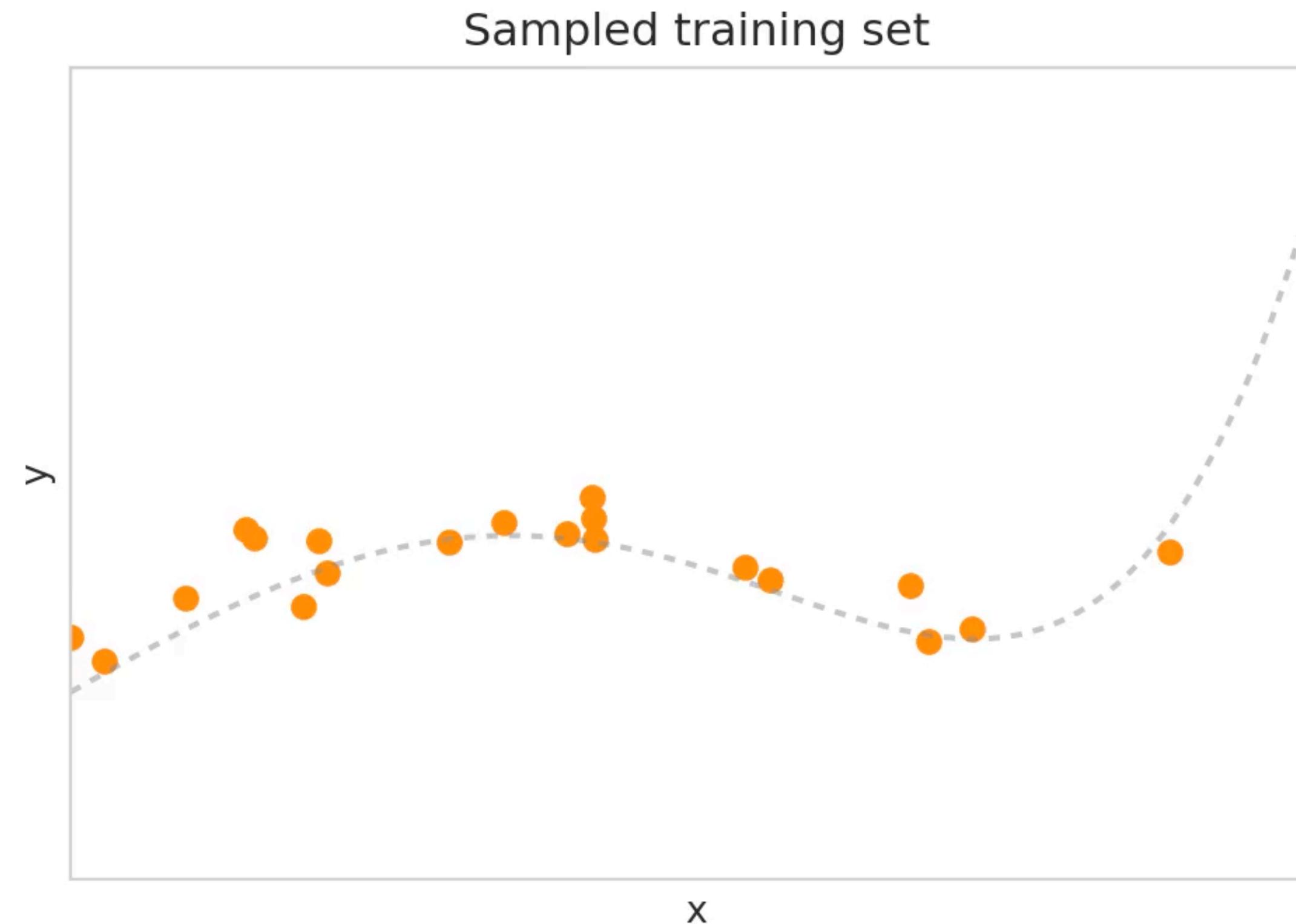
We have observed one particular  $S_{\text{train}}$  but we could have observed several others!

# But there is randomness in the data



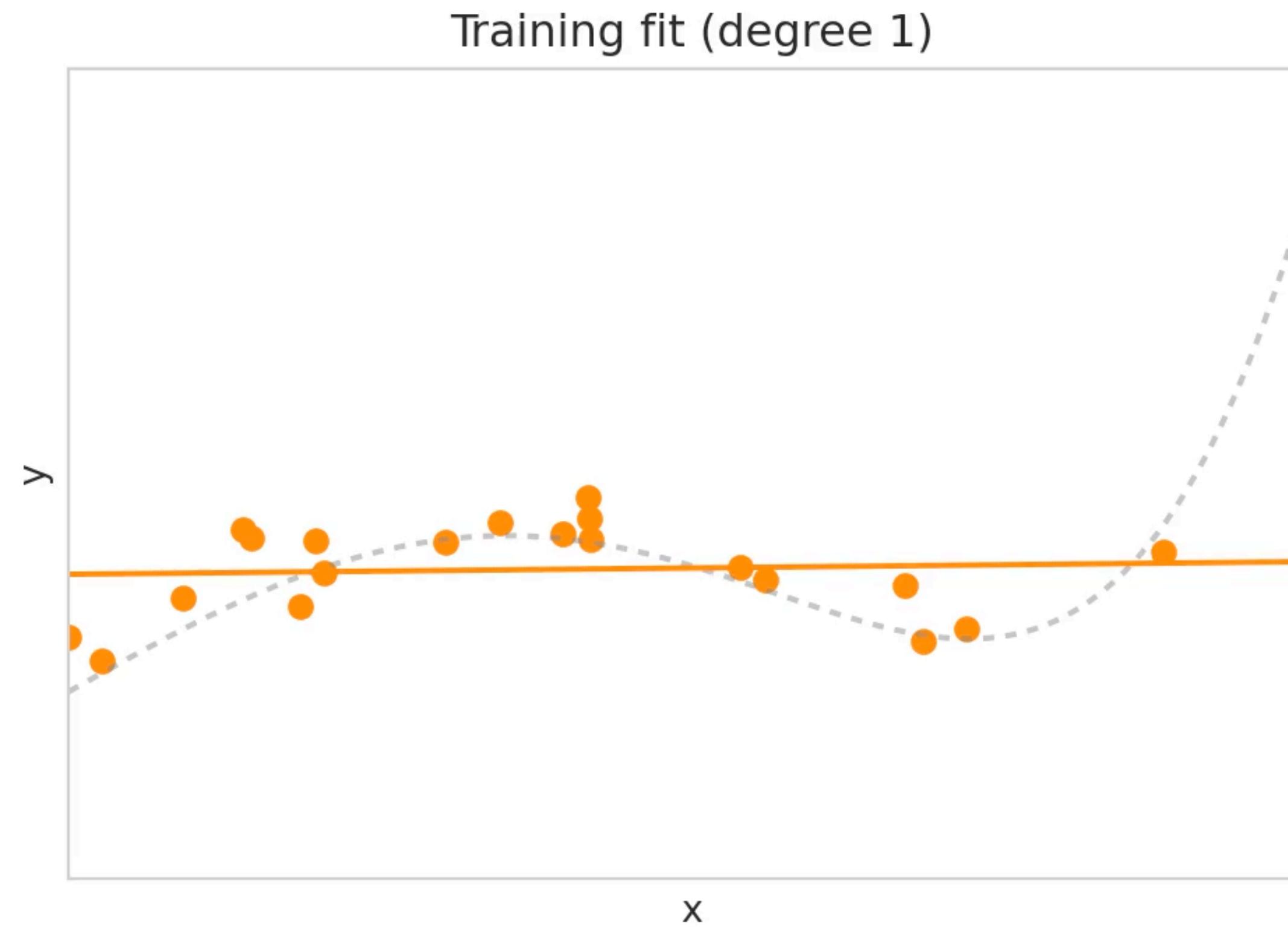
Even if we keep the same  $(x_1, \dots, x_n)$ , we have variability in the observed  $(y_1, \dots, y_n)$

# But there is randomness in the data



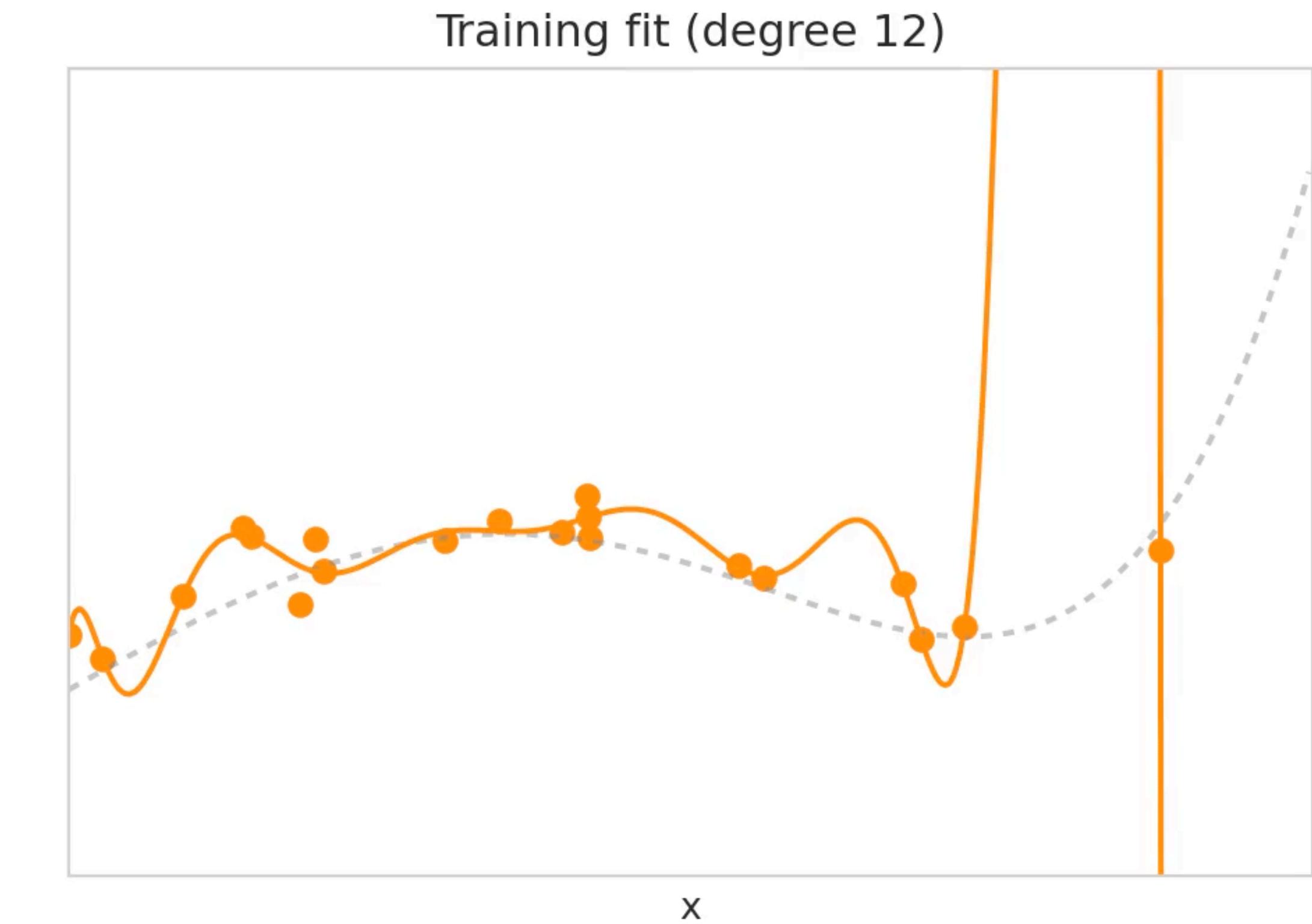
Even if we keep the same  $(x_1, \dots, x_n)$ , we have variability in the observed  $(y_1, \dots, y_n)$

# Thus there is randomness in the predictions



Moving a single observation will cause only a small shift in the position of the line

**Underfitting**

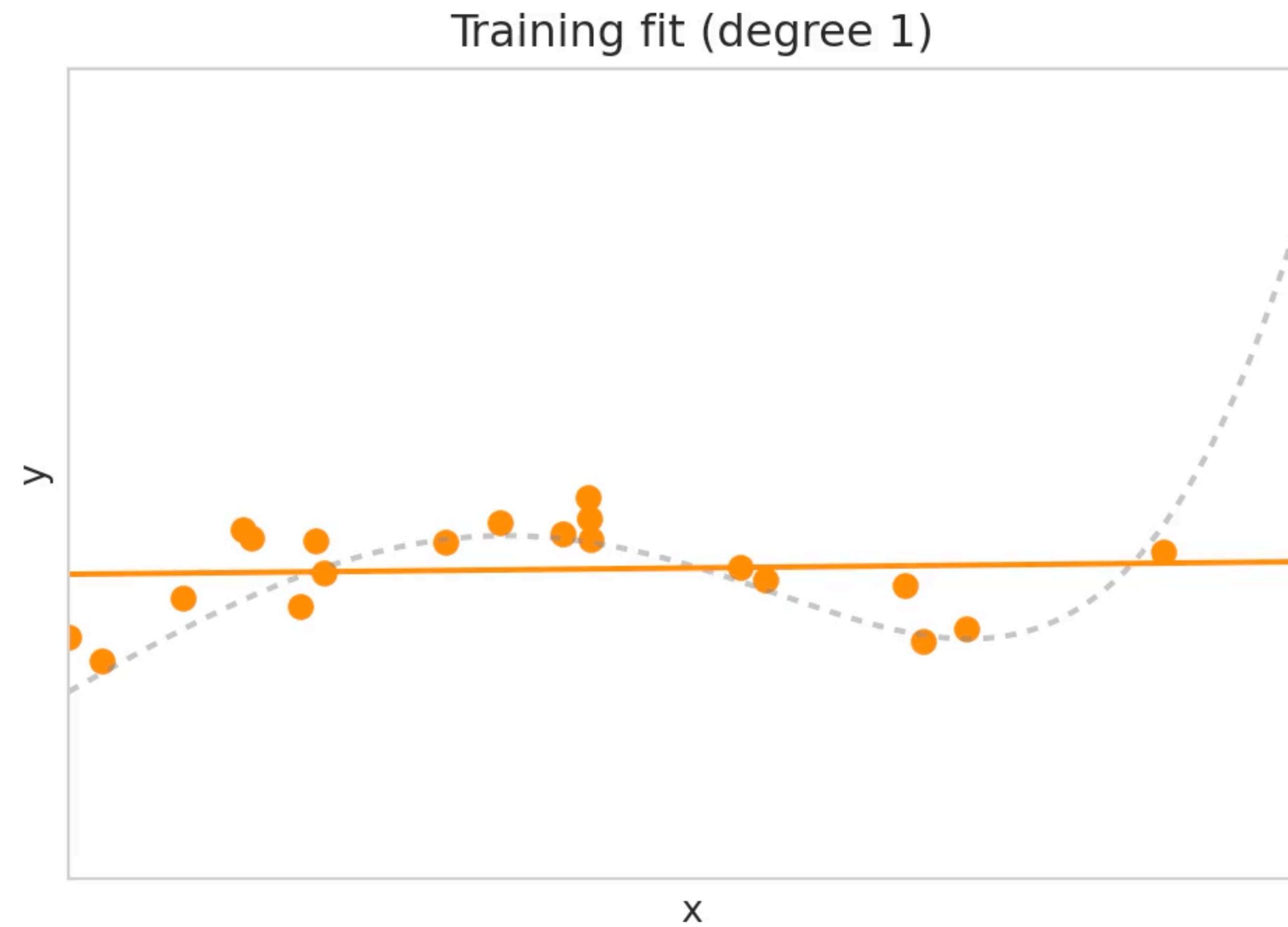


Changing one of the observations may change the prediction considerably

**Overfitting**

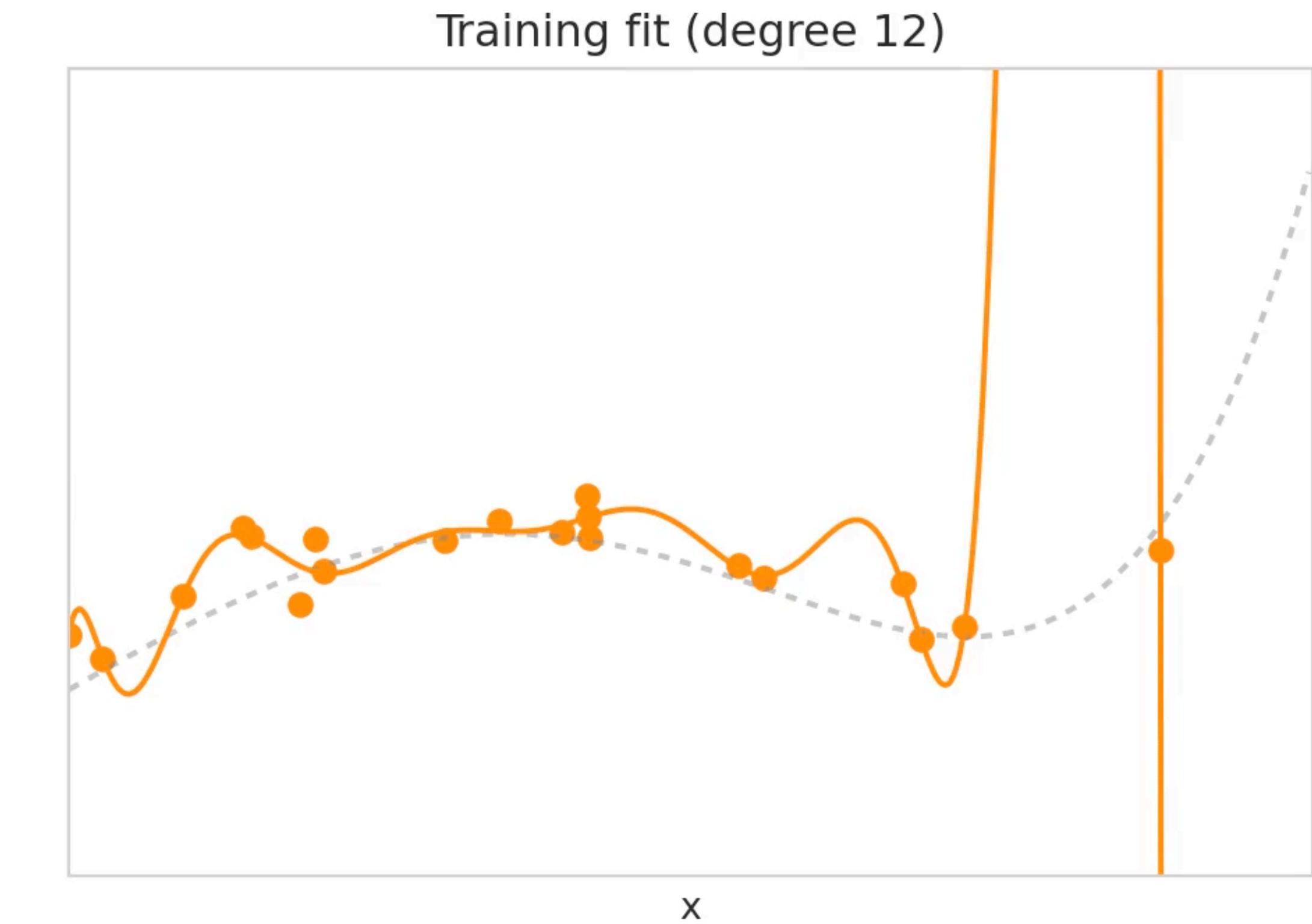
**Simple models are less sensitive**

# Thus there is randomness in the predictions



Moving a single observation will cause only a small shift in the position of the line

**Underfitting**

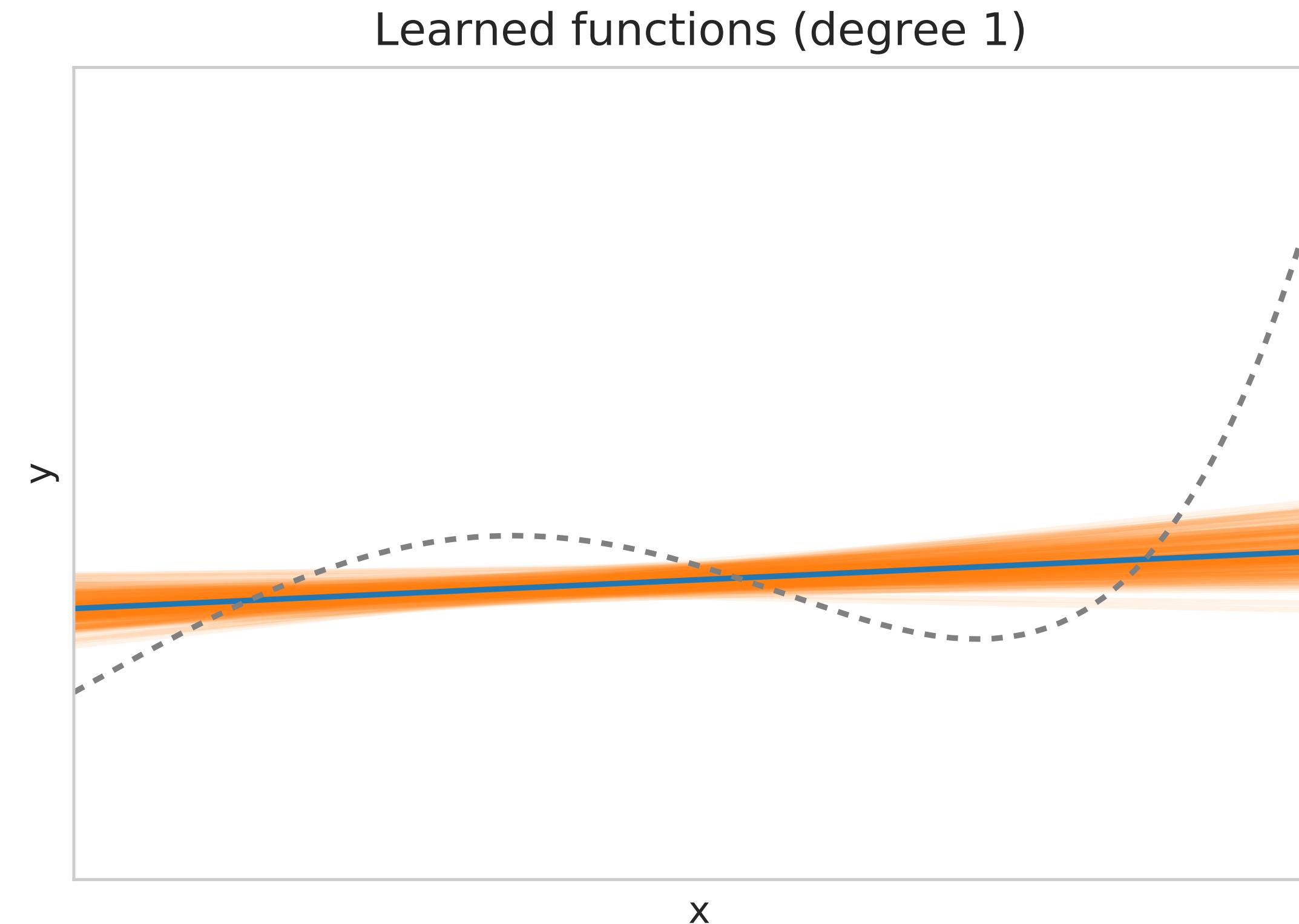


Changing one of the observations may change the prediction considerably

**Overfitting**

**Simple models are less sensitive**

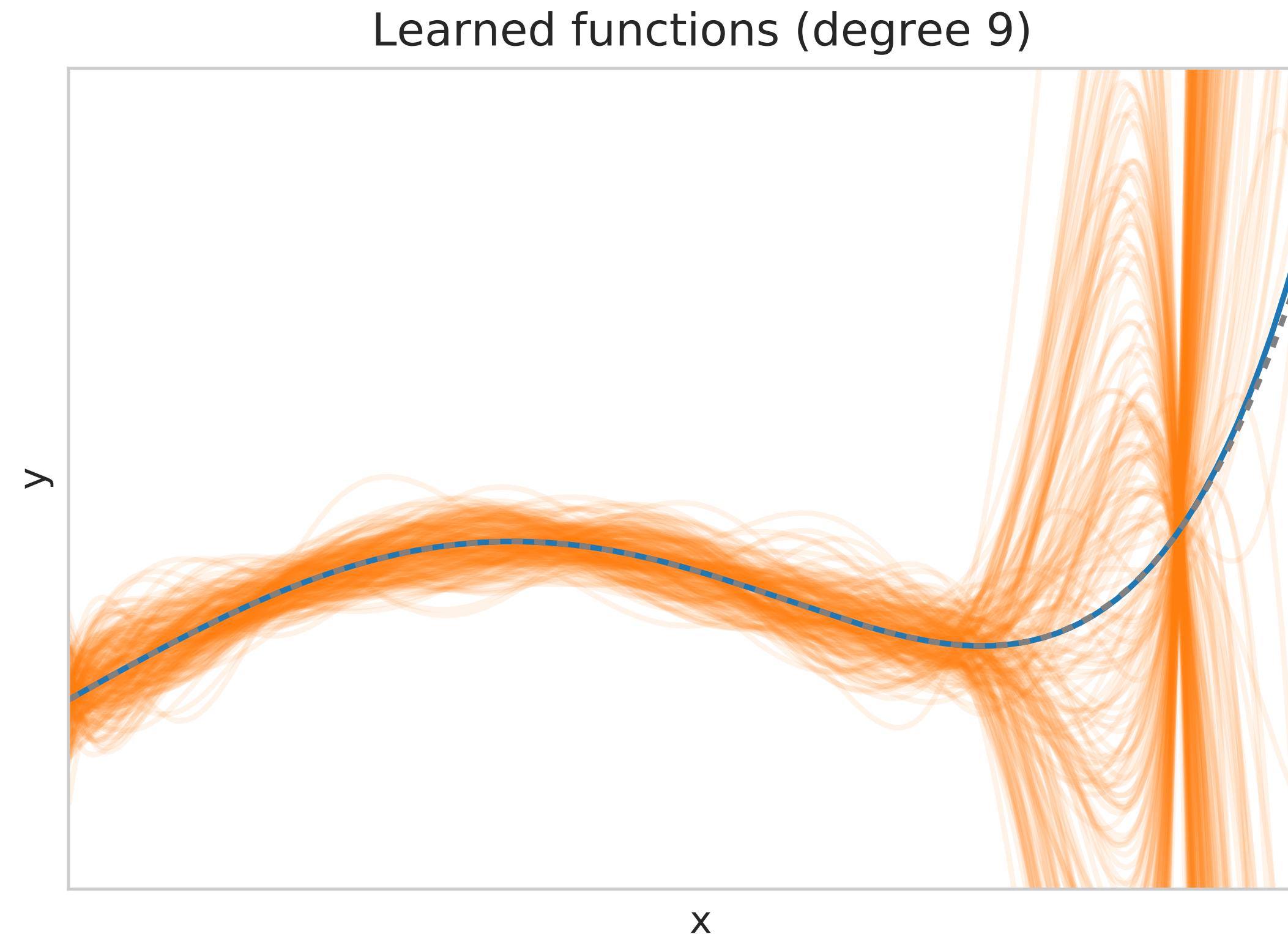
# Simple models have large bias but low variance



The average of the predictions  $f_S$  does not fit the data well: **large bias**

The variance of the predictions  $f_S$  as a function of  $S$  is small: **small variance**

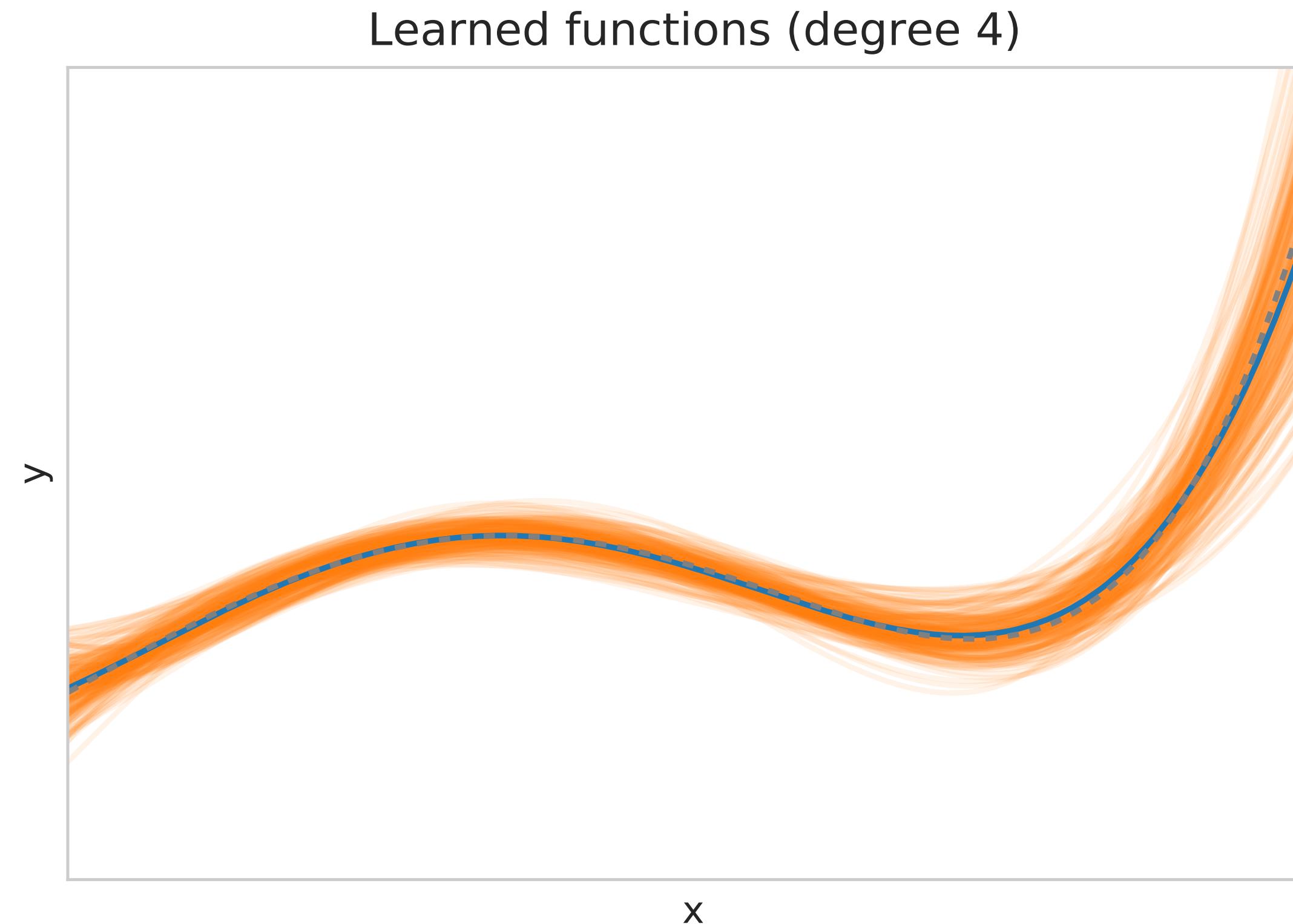
# Complex models have low bias but high variance



The average of the predictions  $f_S$  fits the data well: **small bias**

The variance of the predictions  $f_S$  as a function of S is large: **large variance**

# We need to balance bias & variance correctly



# Data model: output perturbed by some noise

## Output

True model perturbed by some noise

Joint distribution  $(x, y) \sim \mathcal{D}$

## True model

$f$  arbitrary and unknown function

The model is generally **not realizable**,  
i.e.,  $f$  is not in our model class

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

## Input

$$x \sim \mathcal{D}_x$$

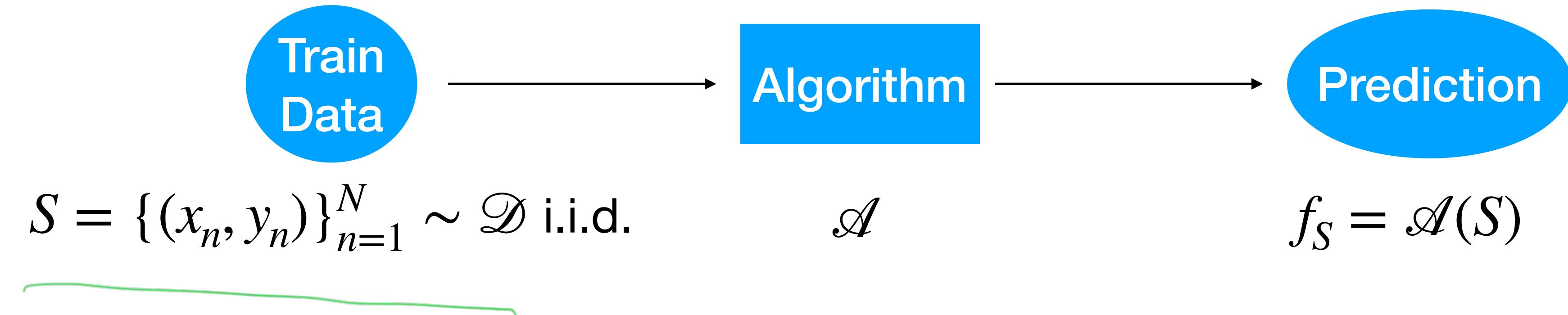
( $\mathcal{D}_x$  fixed but unknown)

## Noise

$\varepsilon \sim \mathcal{D}_\varepsilon$  i.i.d.,  
independent of  $x$   
 $\mathbb{E}[\varepsilon] = 0$

We consider the square loss and will provide a decomposition of the true error

# Error Decomposition



$$p(x, y) = p(x)p(y|x)$$

We are interested in how the **expected error** of  $f_S$ : *fixed  $S$ : constant*

$\mathbb{E}_{(x,y) \sim \mathcal{D}}[(y - f_S(x))^2]$

$\sum_{(x,y)} p(x,y) \dots$

*random  $S$*

$\mathbb{E}_{x \sim \mathcal{D}_x} \mathbb{E}_{y \sim \mathcal{D}_{y|x}} [\dots]$

$\sum_x p(x) \sum_y p(y|x)$

behaves as a **function of the train set  $S$**  and model class complexity

# Error Decomposition



$$S = \{(x_n, y_n)\}_{n=1}^N \sim \mathcal{D} \text{ i.i.d.}$$

$\mathcal{A}$

$$f_S = \mathcal{A}(S)$$

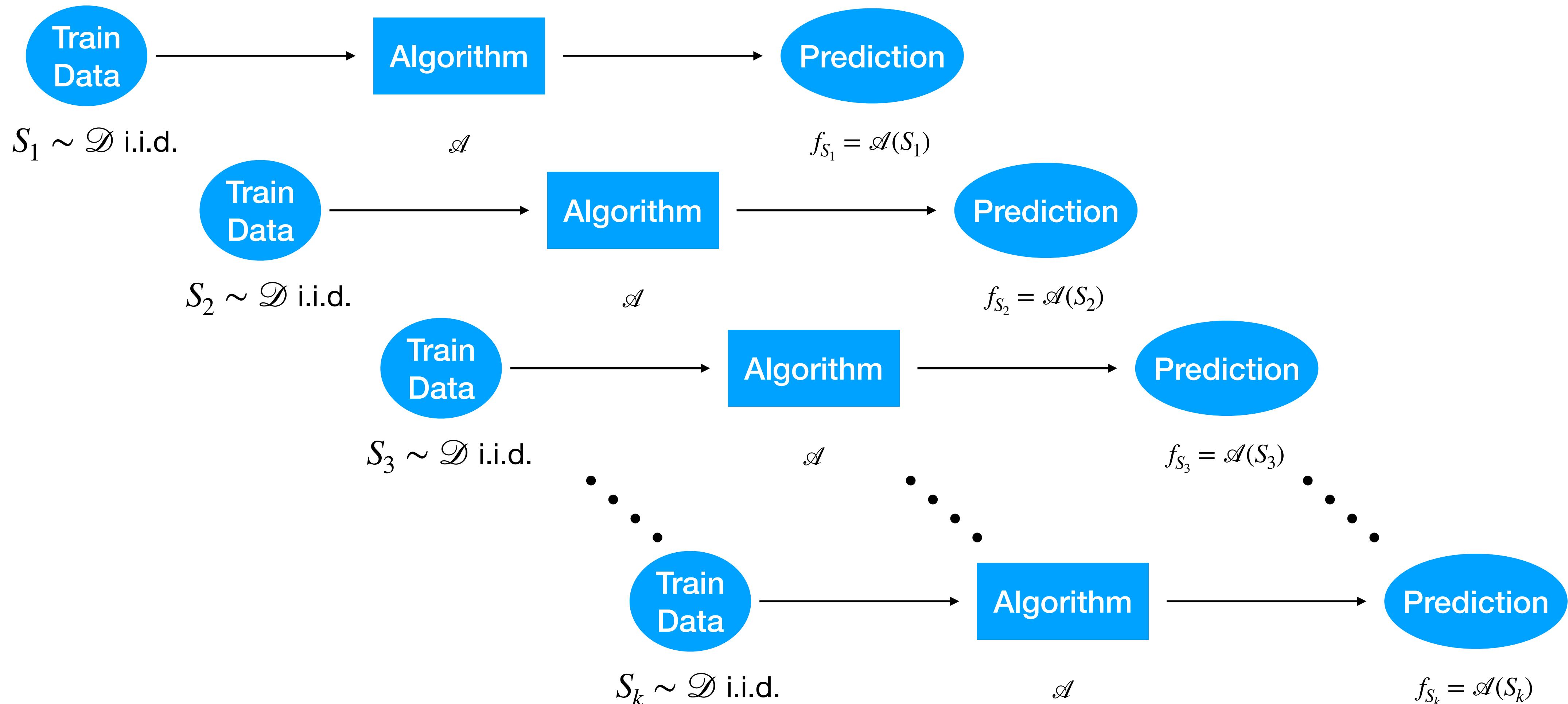
The decomposition will hold true at **every single point**  $x$ . Therefore, to simplify, we consider the expected error of  $f_S$  for a fixed element  $x_0$ :

$$L(f_S) = \mathbb{E}_{\varepsilon \sim \mathcal{D}_\varepsilon} [(f(x_0) + \varepsilon - f_S(x_0))^2]$$

$y_0 = f(x_0) + \varepsilon$

This is a random variable. The randomness comes for the train set  $S$

# We run the experiment many times



We are interested in the **average** and the **variance** of the **predictions** ( $f_{S_1}, \dots, f_{S_k}$ ) over these multiple runs

# A decomposition in three terms

We are interested in the expectation of the true risk over the training set  $S$

$(x, y) \sim \mathcal{D}$

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}}[L(f_S)] &= \mathbb{E}_{S \sim \mathcal{D}} \left[ \mathbb{E}_{\varepsilon \sim \mathcal{D}_\varepsilon} [(f(x_0) + \varepsilon - f_S(x_0))^2] \right] \\ &= \mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon} [(f(x_0) + \varepsilon - f_S(x_0))^2] \end{aligned}$$

We will decompose this quantity in ***three non-negative terms*** and will interpret each of these terms

$$\mathbb{E}_x[c] = c$$

First we expand the square:

$$\mathbb{E}_{S,\varepsilon}[\varepsilon^2] = \mathbb{E}_S \mathbb{E}_{\varepsilon}[\varepsilon^2] = \mathbb{E}_{\varepsilon}[\varepsilon^2]$$

const w.r.t. S

$$\mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_{\varepsilon}}[(f(x_0) + \varepsilon - f_S(x_0))^2] = \mathbb{E}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon^2]$$

$\varepsilon + (f(x_0) - f_S(x_0))$

$$(a+b)^2 = a^2 + 2ab + b^2$$

Q.E.D.

$$\begin{aligned} &+ 2\mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon(f(x_0) - f_S(x_0))] \\ &+ \mathbb{E}_{S \sim \mathcal{D}}[(f(x_0) - f_S(x_0))^2] \end{aligned}$$

Using that  $\mathbb{E}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon] = 0$  and  $\varepsilon \perp\!\!\!\perp S$ :

$$\bullet \mathbb{E}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon^2] = \text{Var}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon] = \mathbb{E}_{\varepsilon}[(\varepsilon - \mathbb{E}\varepsilon)^2] = \mathbb{E}_{\varepsilon}[\varepsilon^2]$$

$$\bullet \mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon(f(x_0) - f_S(x_0))] = \mathbb{E}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon] \cdot \mathbb{E}_{S \sim \mathcal{D}}[f(x_0) - f_S(x_0)] = 0$$

Therefore

$$\mathbb{E}[X \cdot Y] = (\mathbb{E}X) \cdot (\mathbb{E}Y) \iff X \perp\!\!\!\perp Y$$

$$\boxed{\mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_{\varepsilon}}[(f(x_0) + \varepsilon - f_S(x_0))^2] = \text{Var}_{\varepsilon \sim \mathcal{D}_{\varepsilon}}[\varepsilon] + \mathbb{E}_{S \sim \mathcal{D}}[(f(x_0) - f_S(x_0))^2]}$$

$$(a+b)^2 = a^2 + b^2 + 2ab$$

Trick: we add and subtract the constant term  $\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]$ , where  $S'$  is a second training set independent from  $S$

$$= \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)]$$

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}}[(f(x_0) - f_S(x_0))^2] &= \mathbb{E}_{S \sim \mathcal{D}}\left[\left[f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]\right] + \left[\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0)\right]^2\right] \\ &= \mathbb{E}_{S \sim \mathcal{D}}\left[\underbrace{(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)])^2}_{a^2} + \underbrace{(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0))^2}_{b^2}\right] \\ &\quad + 2(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)])(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0)) \end{aligned}$$

Cross-term:

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}}\left[\left(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]\right) \cdot \left(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0)\right)\right] \\ &= (f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]) \cdot \mathbb{E}_{S \sim \mathcal{D}}\left[\left(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0)\right)\right] \\ &= (f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]) \cdot (\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)]) = 0. \end{aligned}$$

$$\boxed{\mathbb{E}_{S \sim \mathcal{D}}[(f(x_0) - f_S(x_0))^2] = (f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)])^2 + \mathbb{E}_{S \sim \mathcal{D}}[(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0))^2]}$$

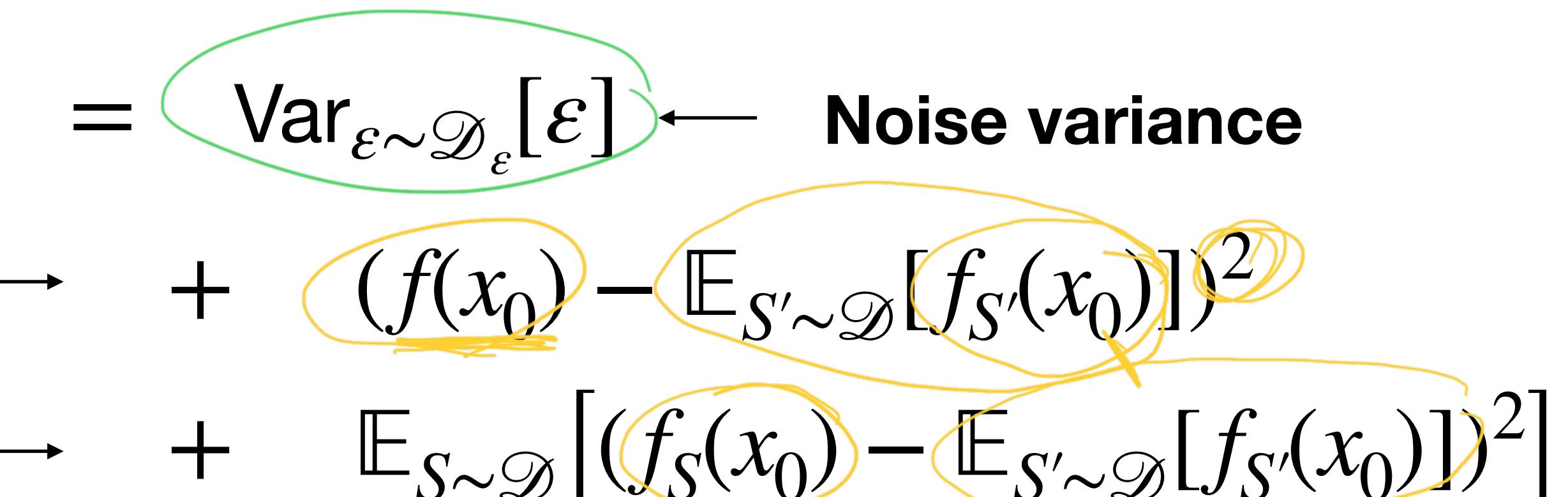
# Bias-Variance Decomposition

We obtain the following decomposition into three positive terms:

$$\mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon} [(f(x_0) + \varepsilon - f_S(x_0))^2] = \text{Var}_{\varepsilon \sim \mathcal{D}_\varepsilon} [\varepsilon] \leftarrow \textbf{Noise variance}$$

**Bias<sup>2</sup>** → +  $(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}} [f_{S'}(x_0)])^2$

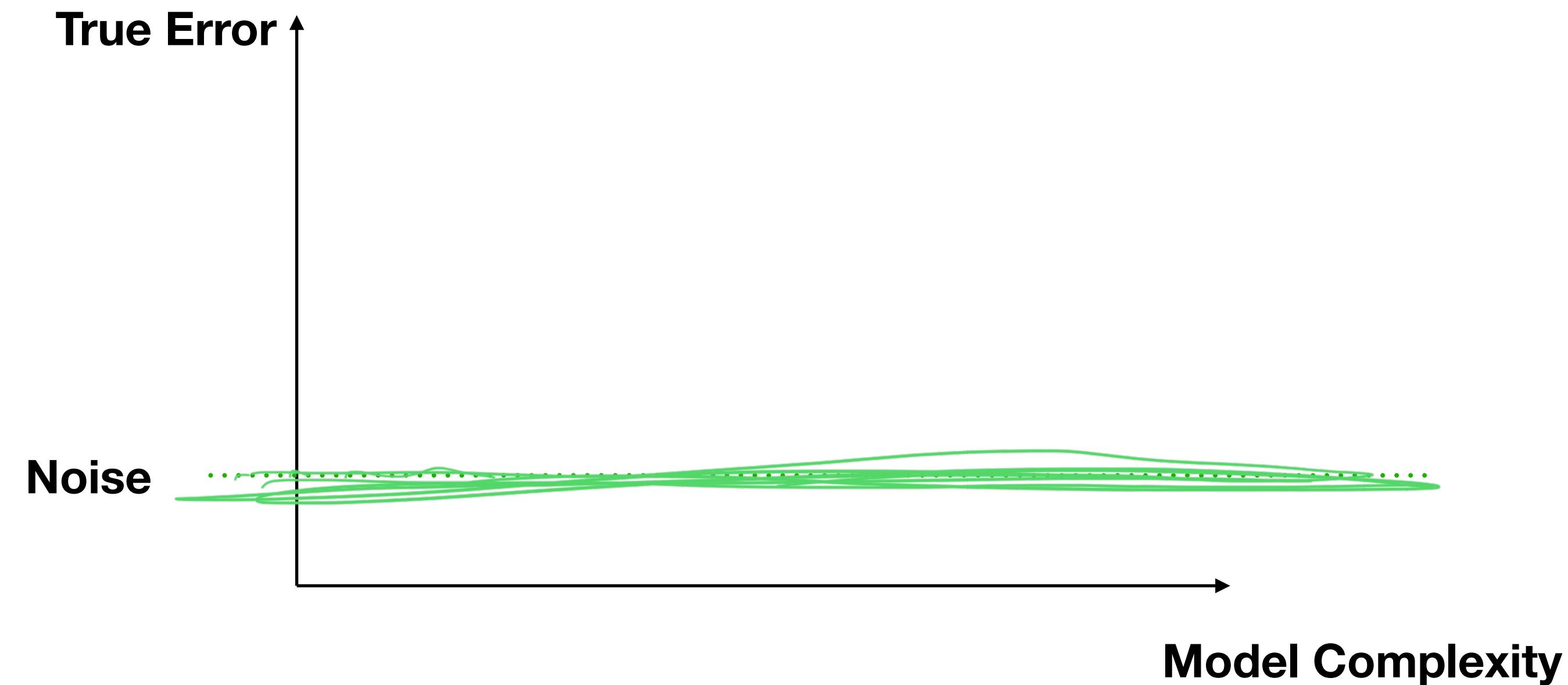
**Variance** → +  $\mathbb{E}_{S \sim \mathcal{D}} [(f_S(x_0) - \mathbb{E}_{S' \sim \mathcal{D}} [f_{S'}(x_0)])^2]$



each of which always provides a lower bound of the true error

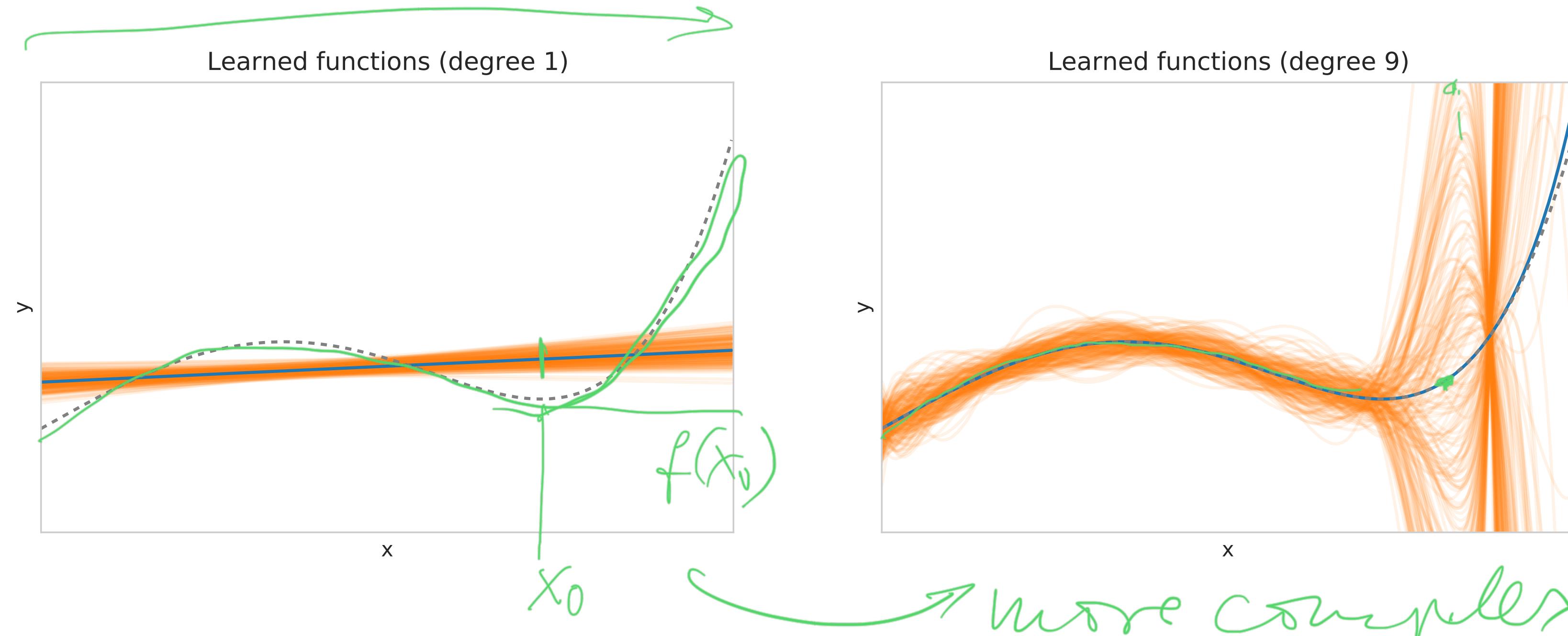
- To minimize the true error, we must choose a method that achieves **low bias and low variance** simultaneously

# Noise: a strict lower bound on the achievable error



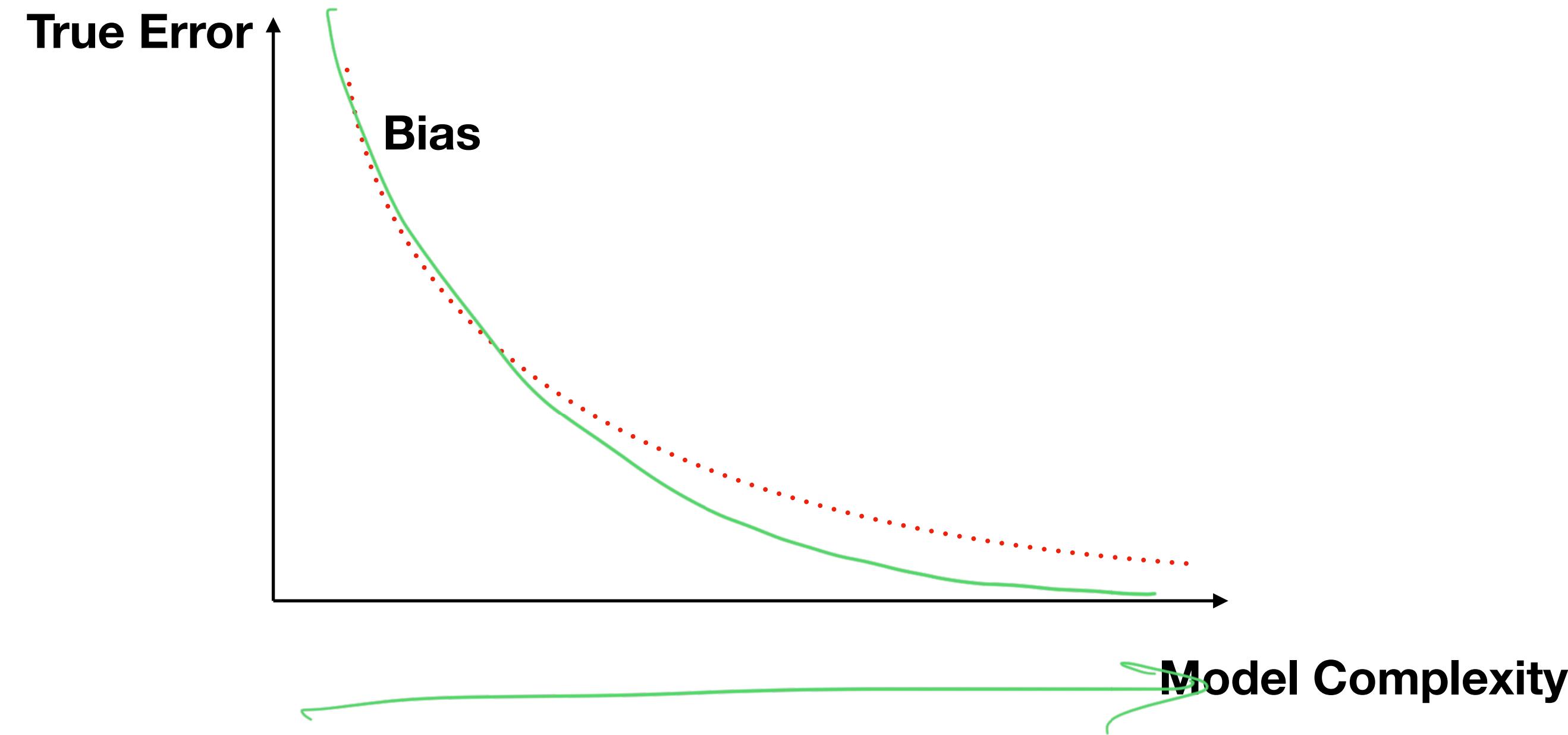
- It is not possible to go below the noise level
- Even if we know the true model  $f$ , we still suffer from the noise:  $L(f) = \mathbb{E}[\varepsilon^2]$
- It is not possible to predict the noise from the data since they are independent

$$\text{Bias: } (f(x_0) - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)])^2$$



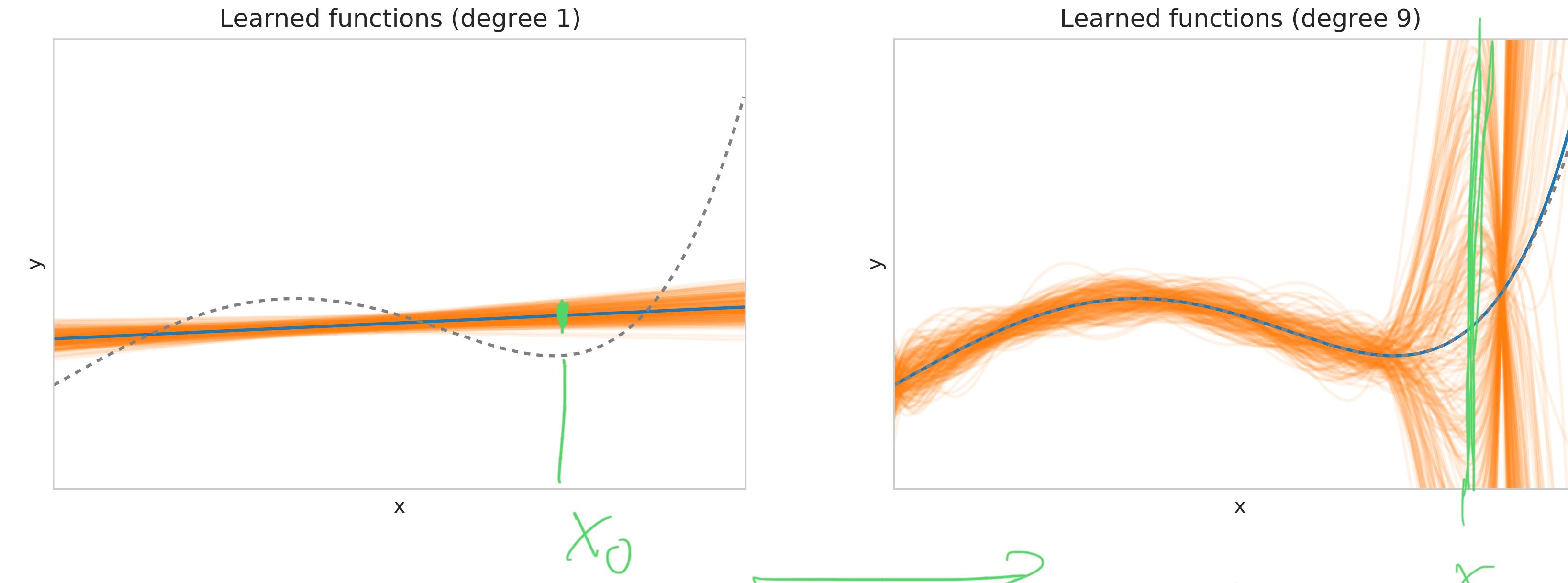
- Squared of the difference between the actual value  $f(x_0)$  and the expected prediction
- It measures how far off in general the models' predictions are from the correct value
- If model **complexity is low**, **bias is typically high**
- If model **complexity is high**, **bias is typically low**

$$\text{Bias: } (f(x_0) - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)])^2$$



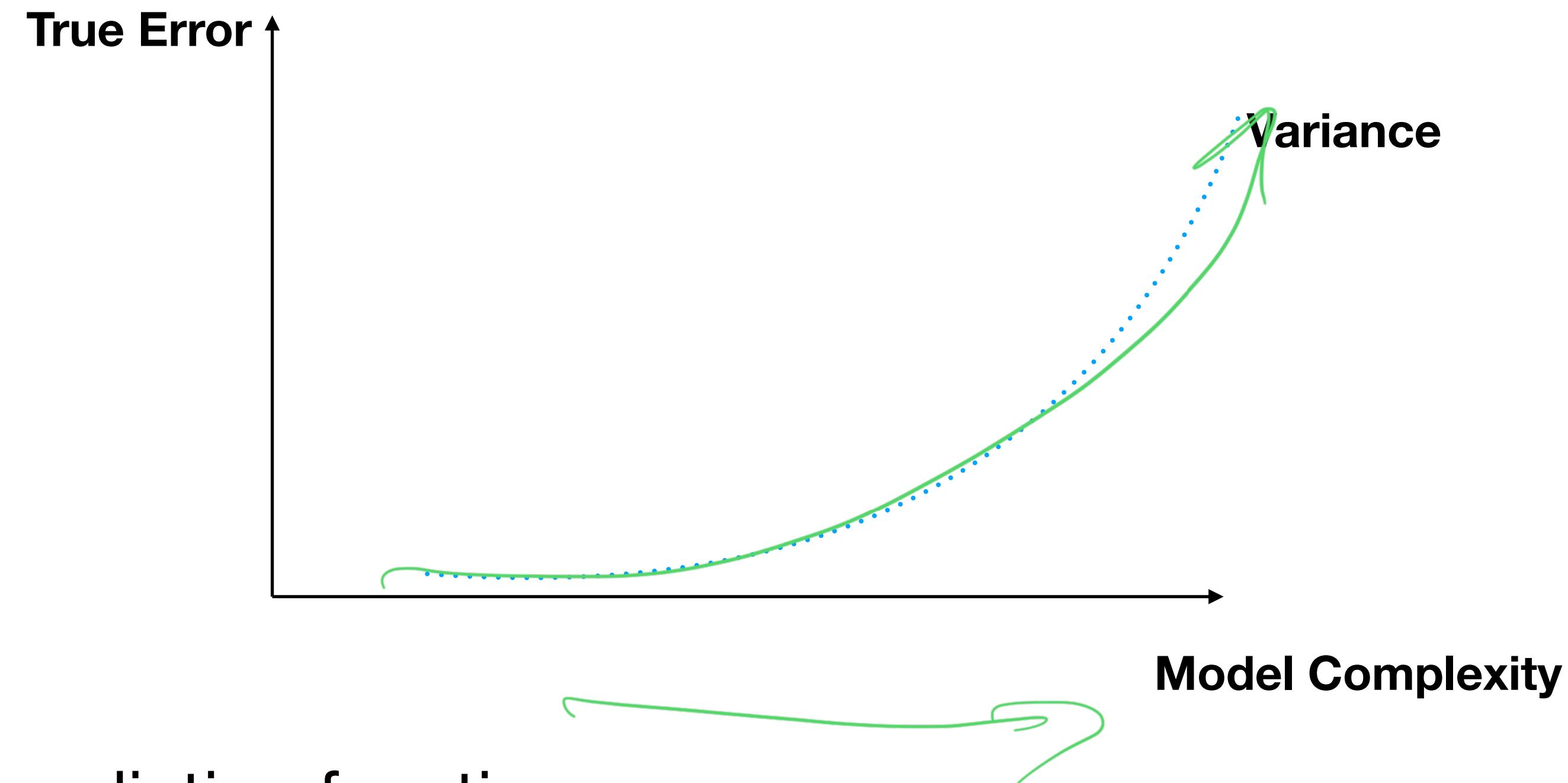
- Squared of the difference between the actual value  $f(x_0)$  and the expected prediction
- It measures how far off in general the models' predictions are from the correct value
- If model **complexity is low**, **bias is typically high**
- If model **complexity is high**, **bias is typically low**

$$\text{Variance: } \mathbb{E}_{S \sim \mathcal{D}}[(f_S(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)])^2]$$



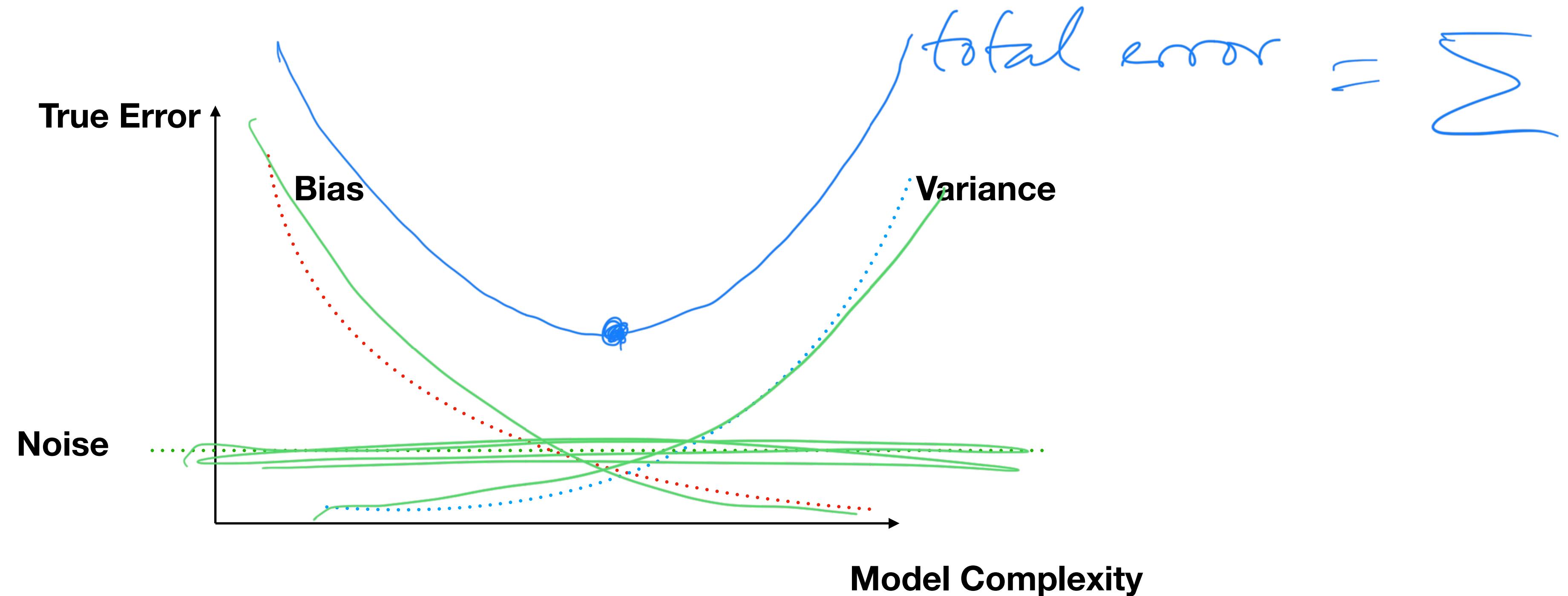
- Variance of the prediction function
- It measures the variability of predictions at a given point across different training set realizations
- If we consider complex models, small variations in the training set can lead to significant changes in the predictions

$$\text{Variance: } \mathbb{E}_{S \sim \mathcal{D}} [(f_S(x_0) - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)])^2]$$



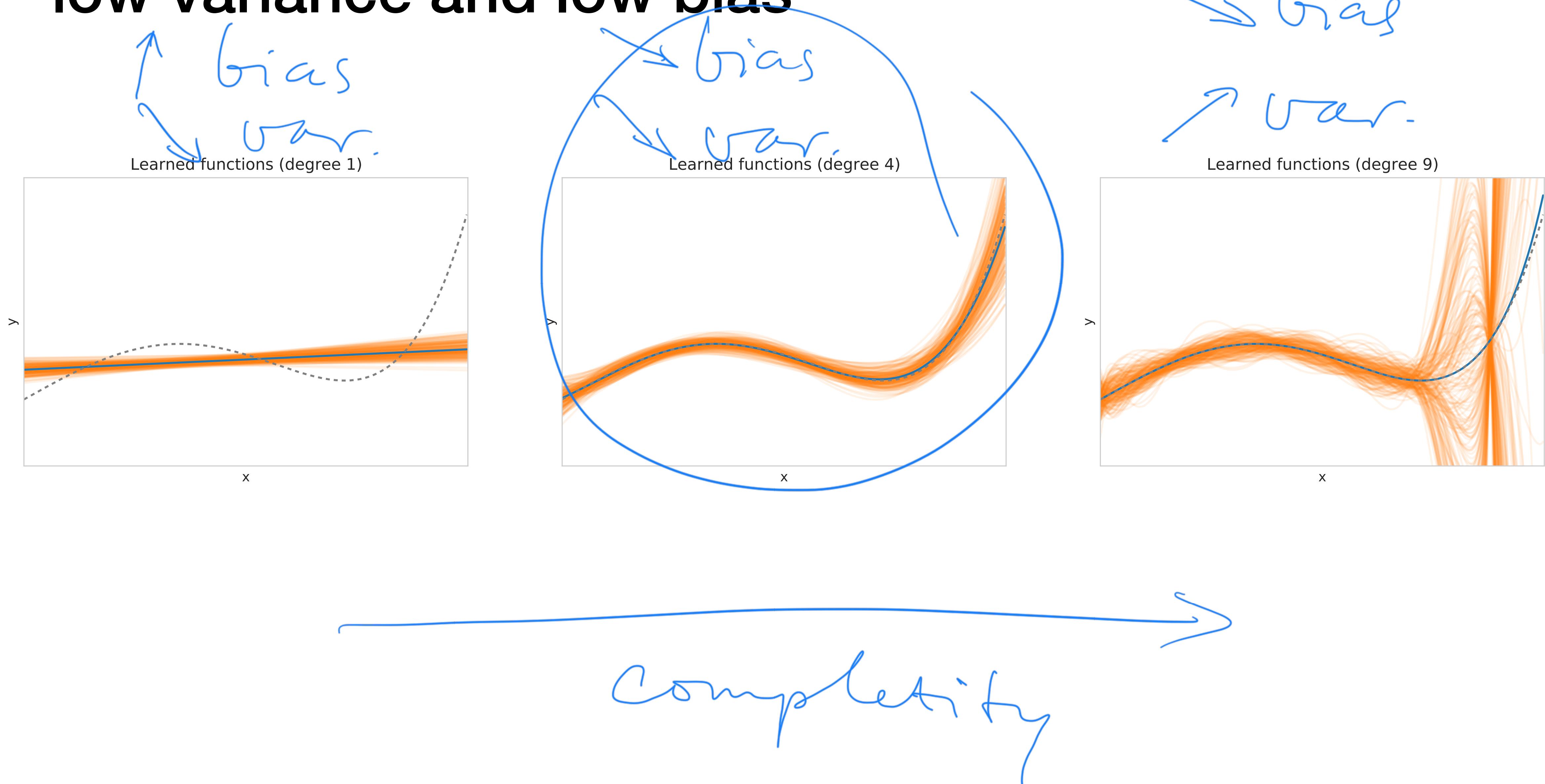
- Variance of the prediction function
- It measures the variability of predictions at a given point across different training set realizations
- If we consider complex models, small variations in the training set can lead to significant changes in the predictions

# Bias Variance tradeoff and U-shape curve



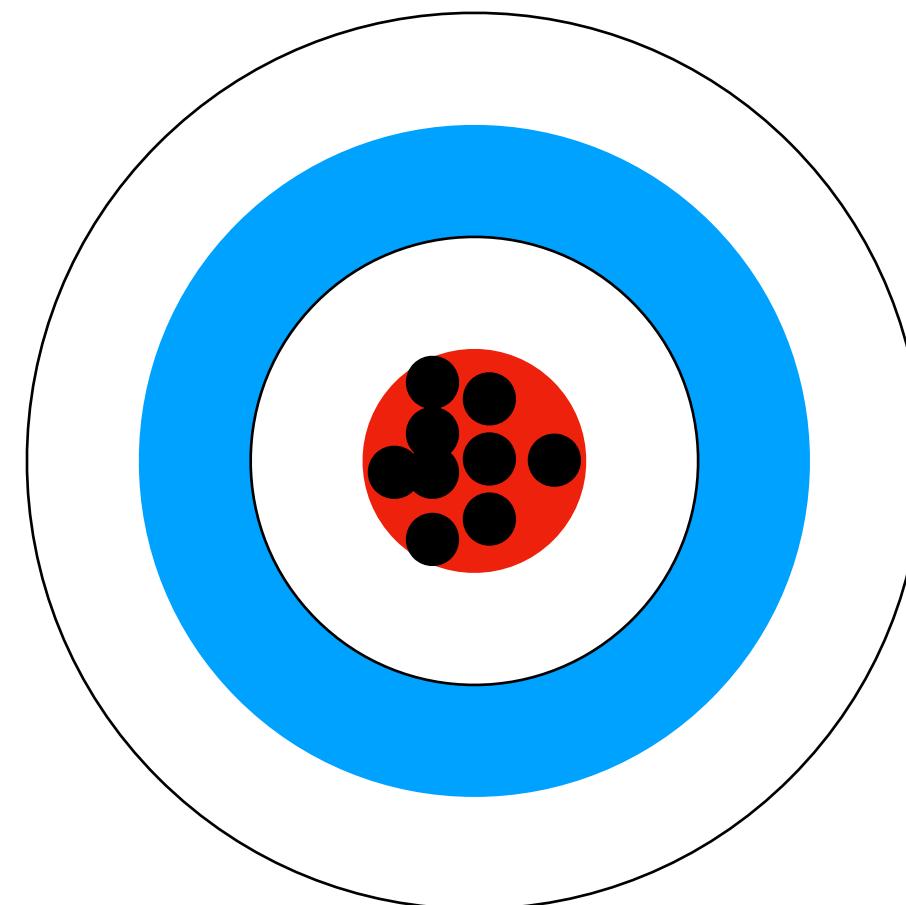
- If model complexity is too low, approximation will be poor (underfitting)
- If model complexity is too high, it may cause issues with variance (overfitting)
  - ➡ This phenomenon is known as the bias-variance tradeoff

# Challenge: Identify a method that ensures both low variance and low bias

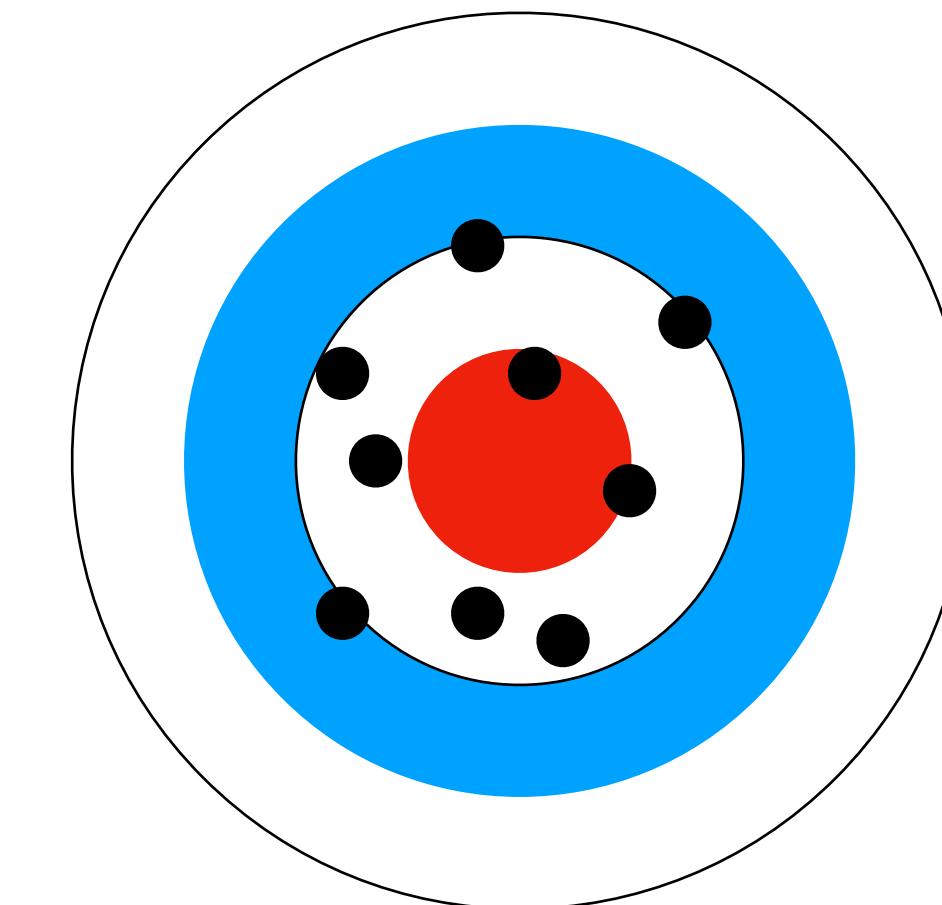


**Low Bias**

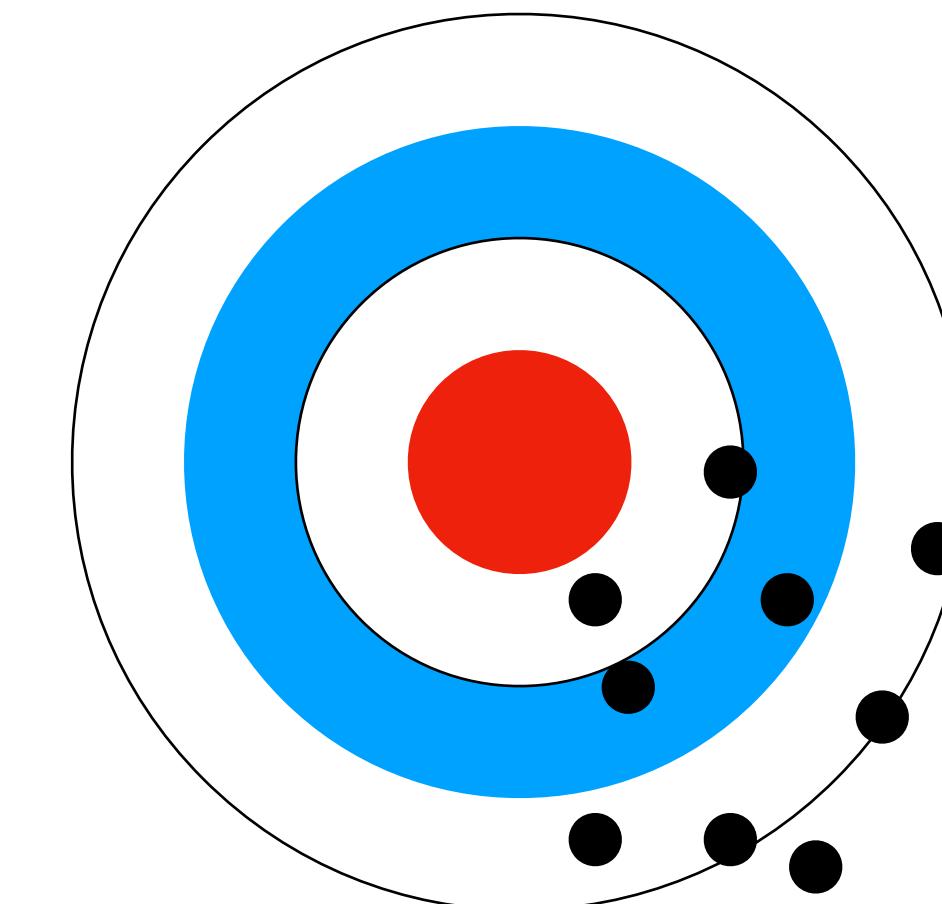
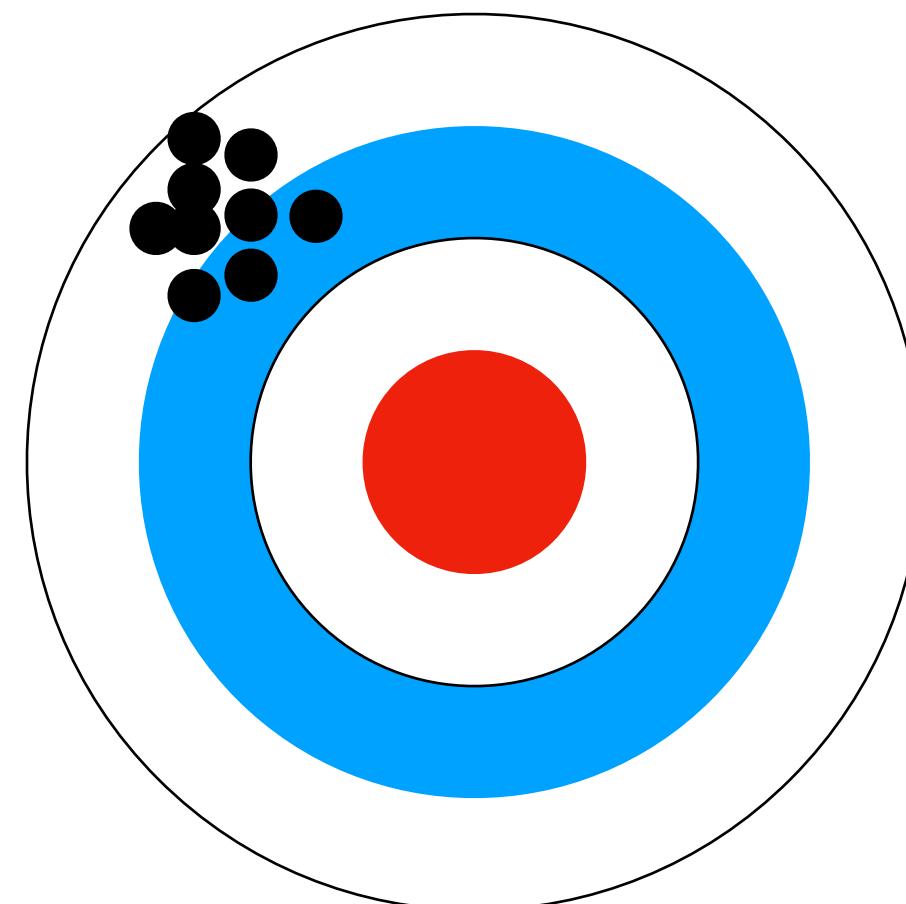
**Low Variance**



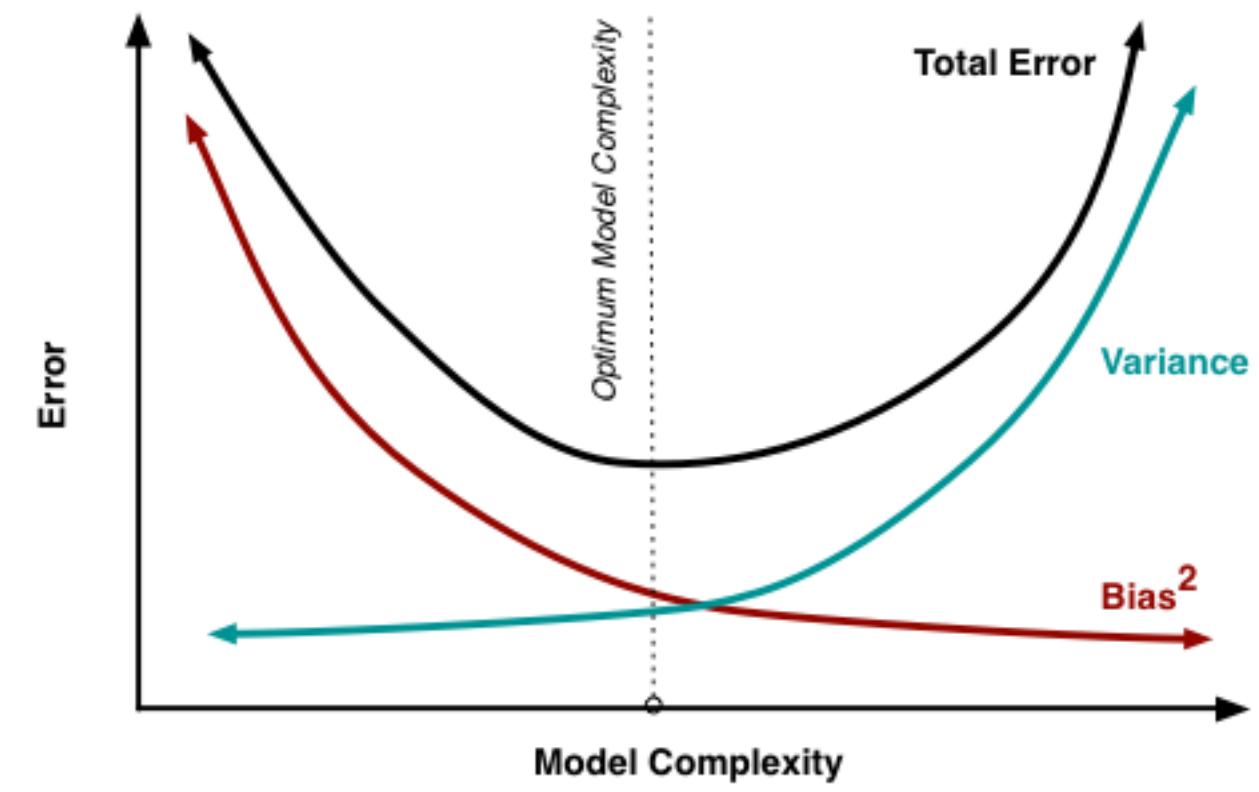
**High Variance**



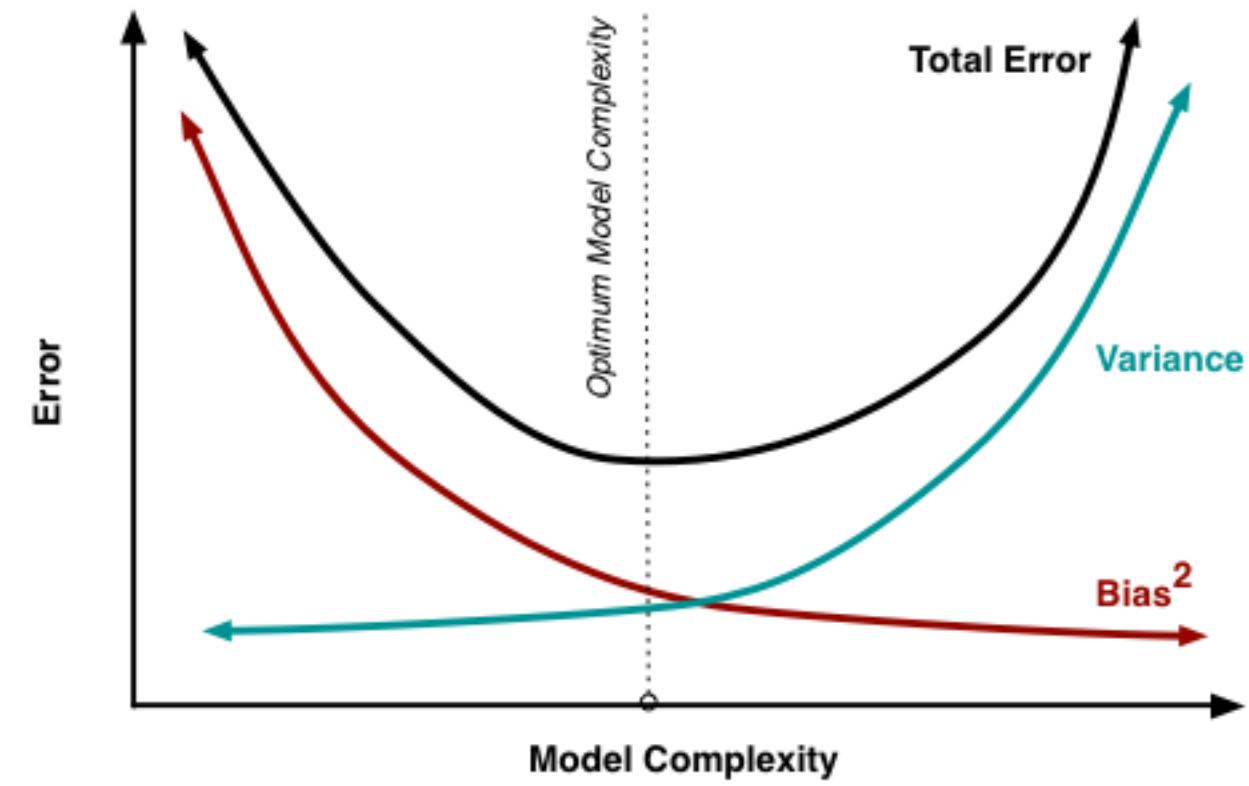
**High Bias**



# In practice...



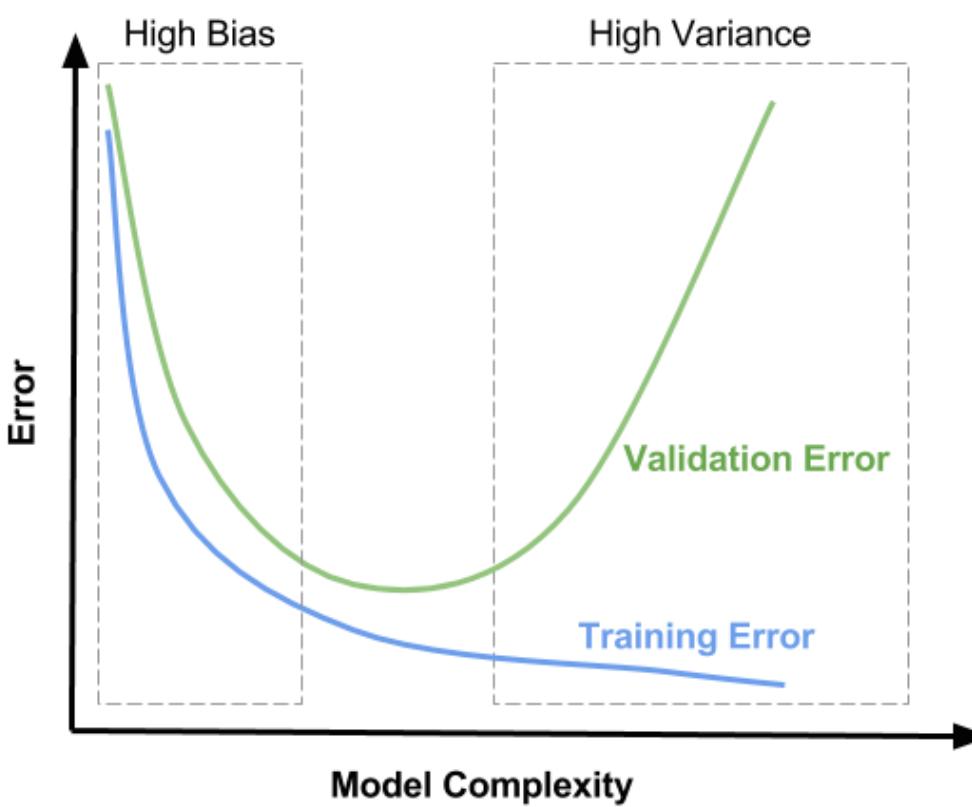
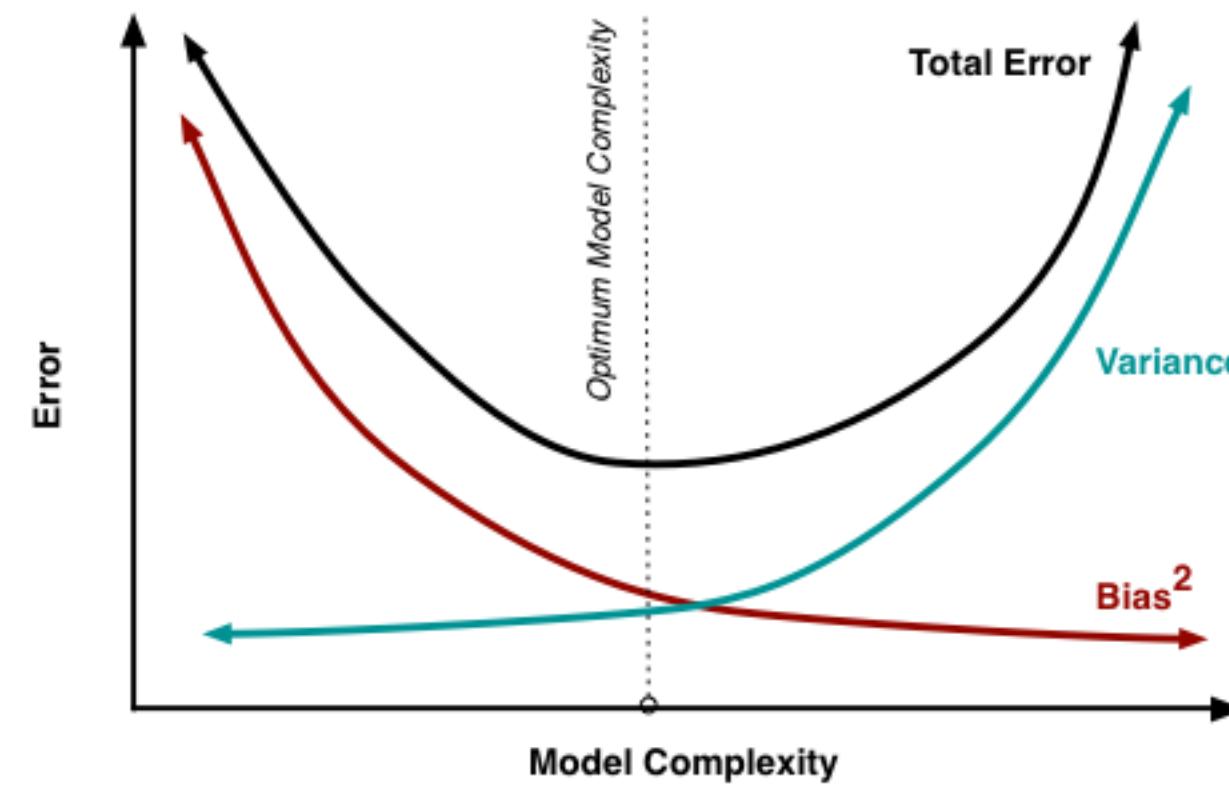
# In practice...



Don't know bias!

If model is trained  
only once per  
complexity level:  
don't have  
variance!

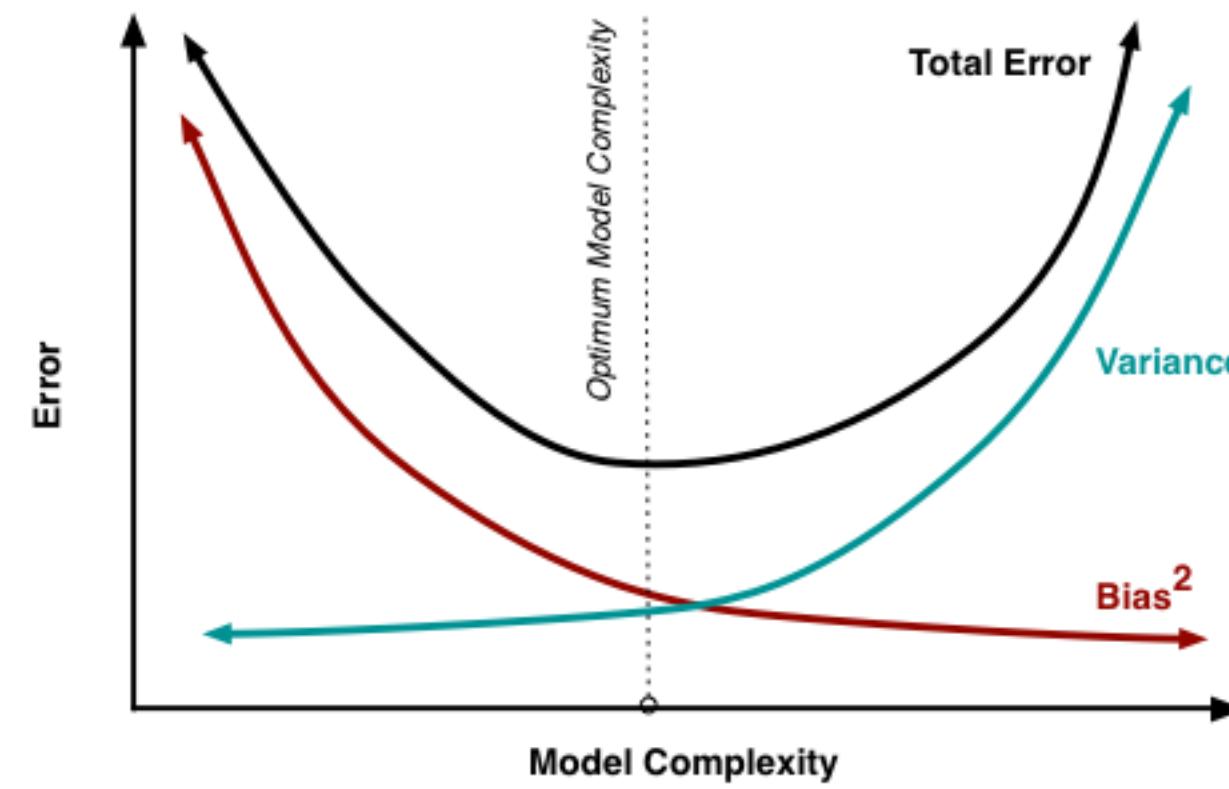
# In practice...



Don't know bias!

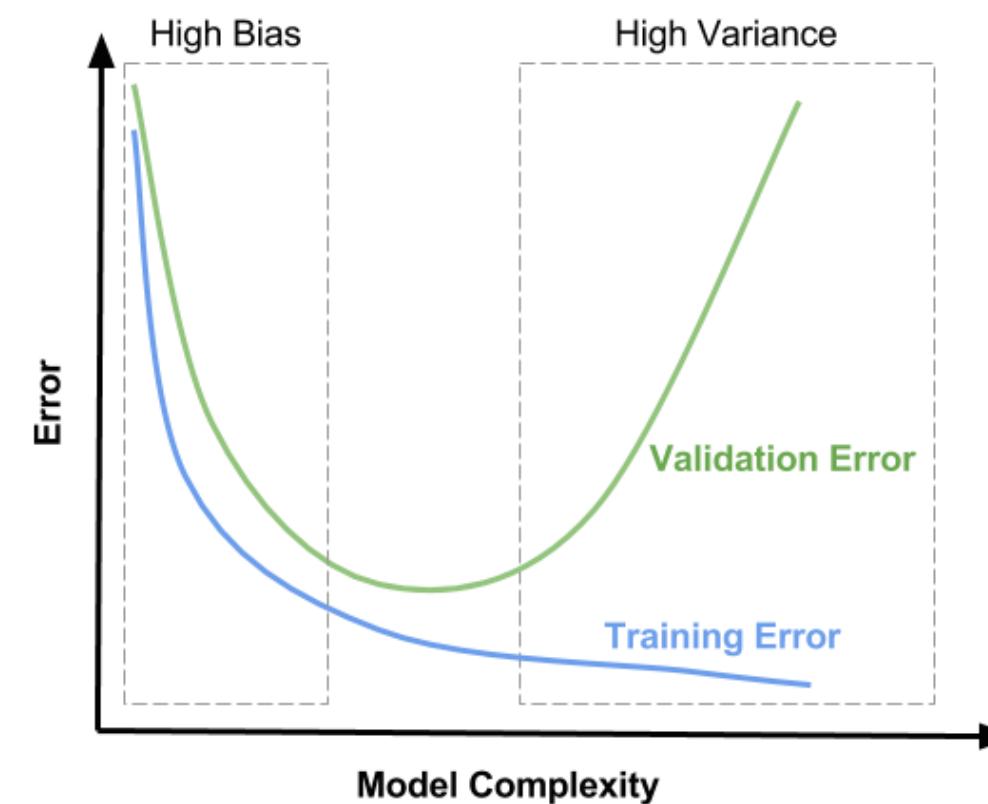
If model is trained  
only once per  
complexity level:  
don't have  
variance!

# In practice...



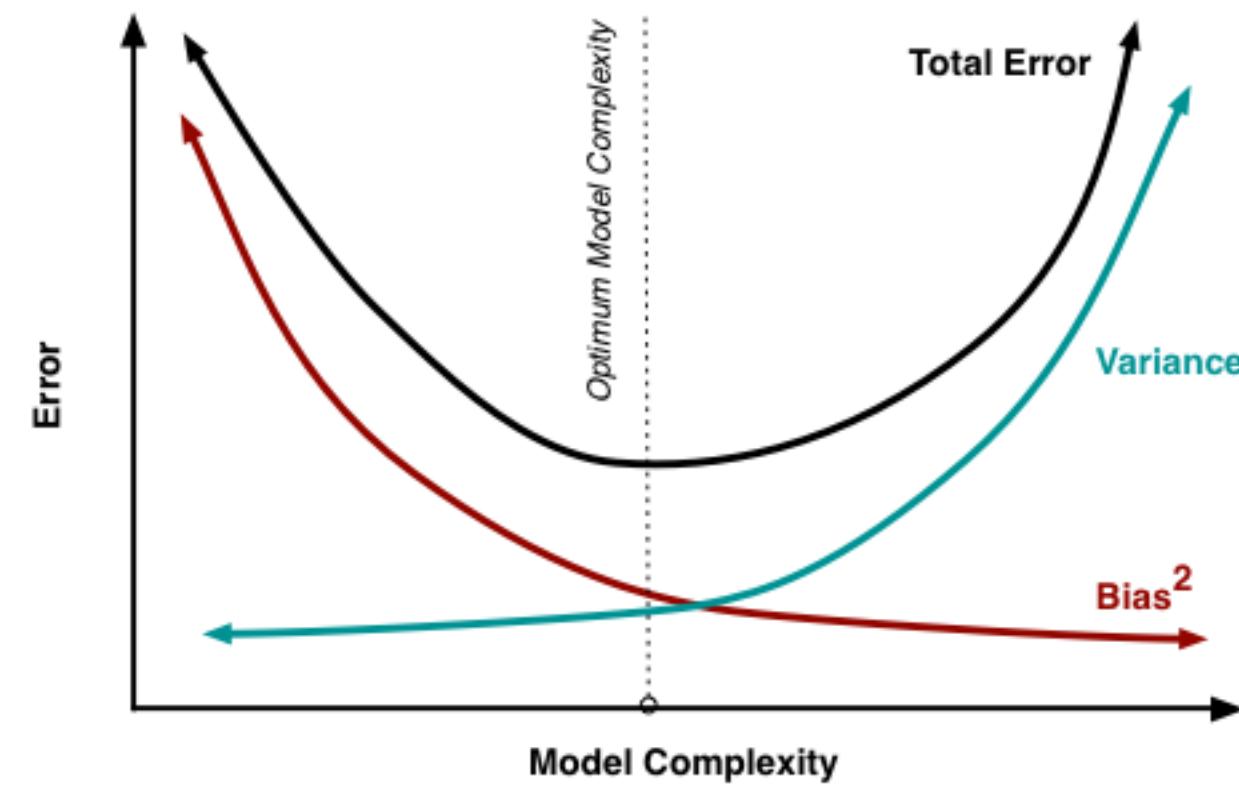
Don't know bias!

If model is trained  
only once per  
complexity level:  
don't have  
variance!



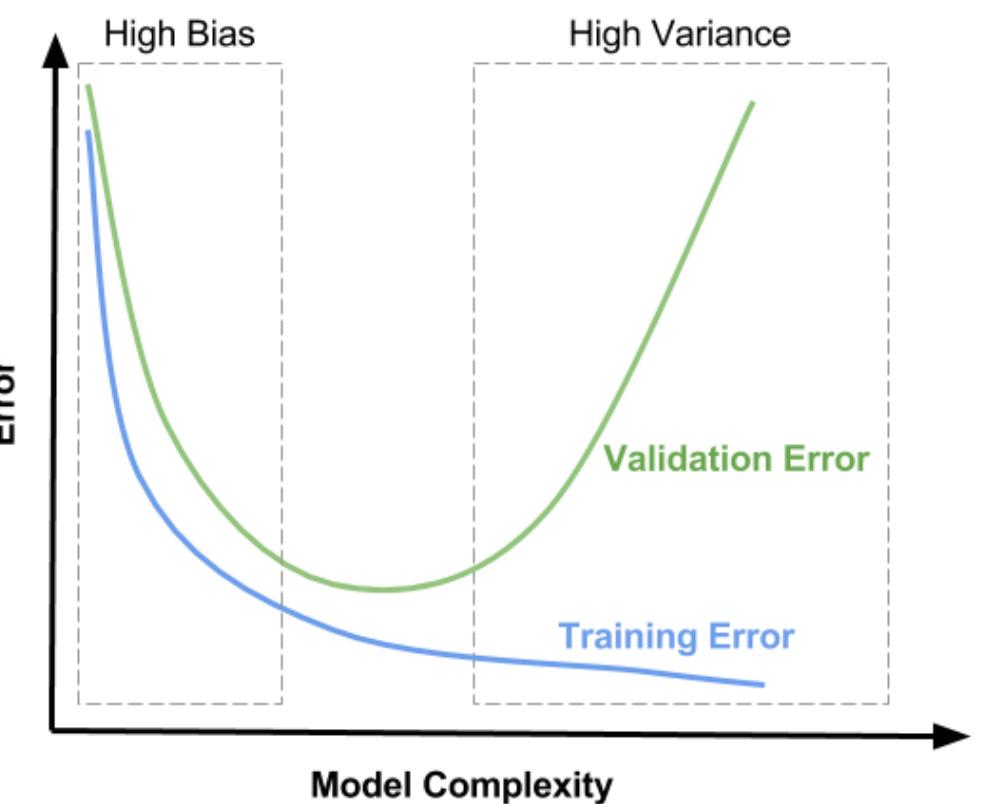
Can't draw curve  
before sweeping  
over complexity  
levels  
(hyperparameter  
values)!

# In practice...

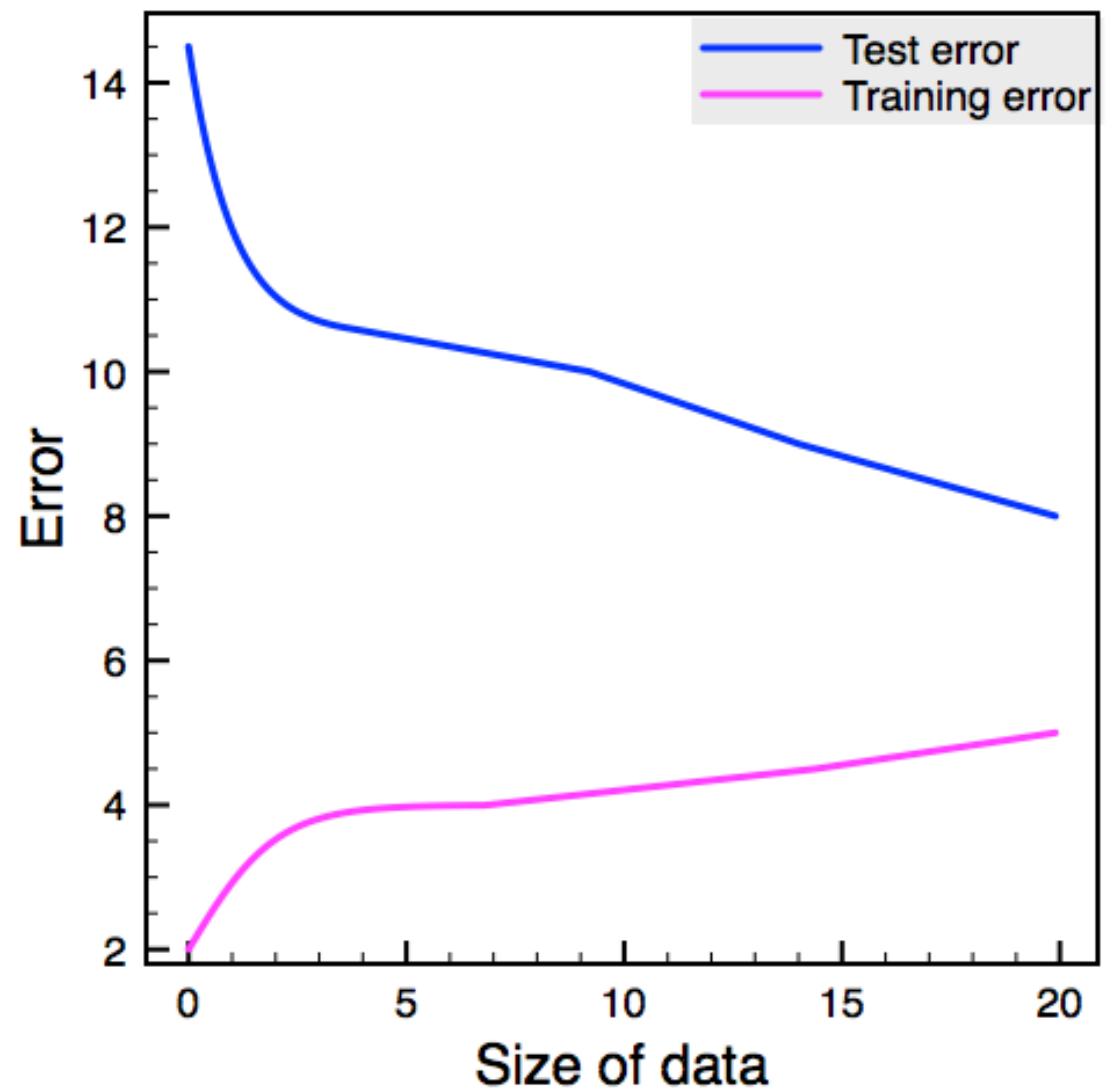


Don't know bias!

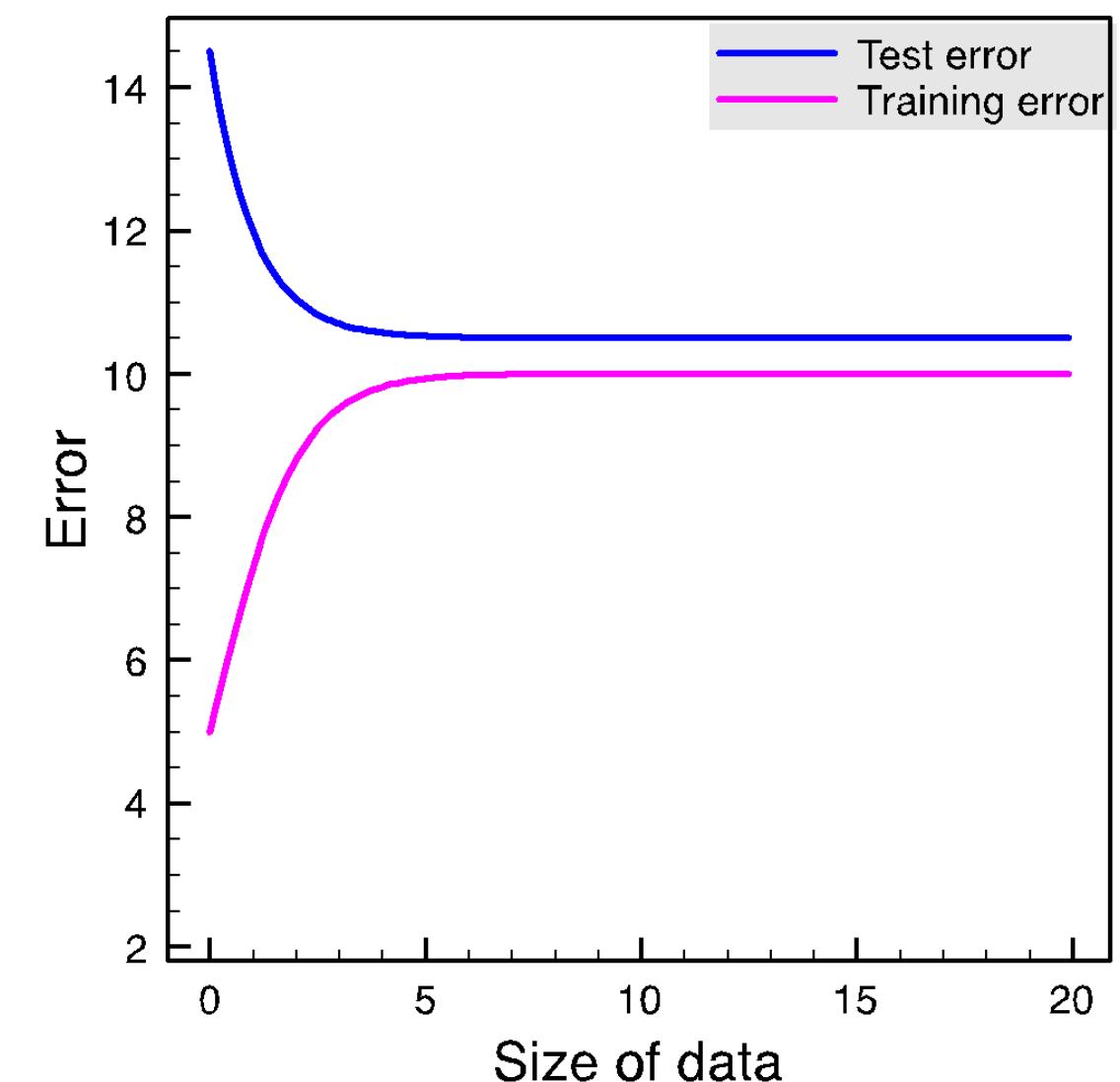
If model is trained  
only once per  
complexity level:  
don't have  
variance!



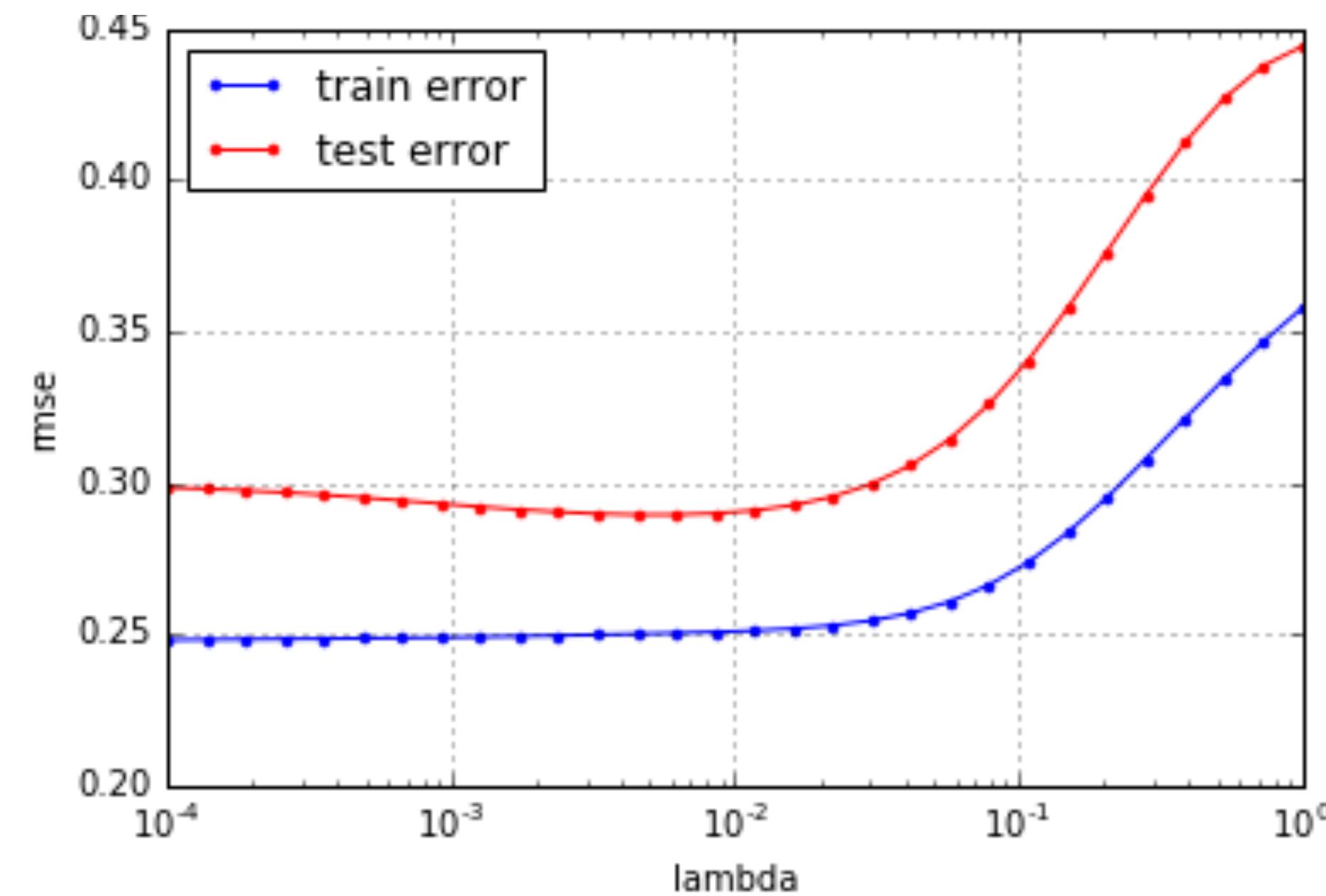
Can't draw curve  
before sweeping  
over complexity  
levels  
(hyperparameter  
values)!



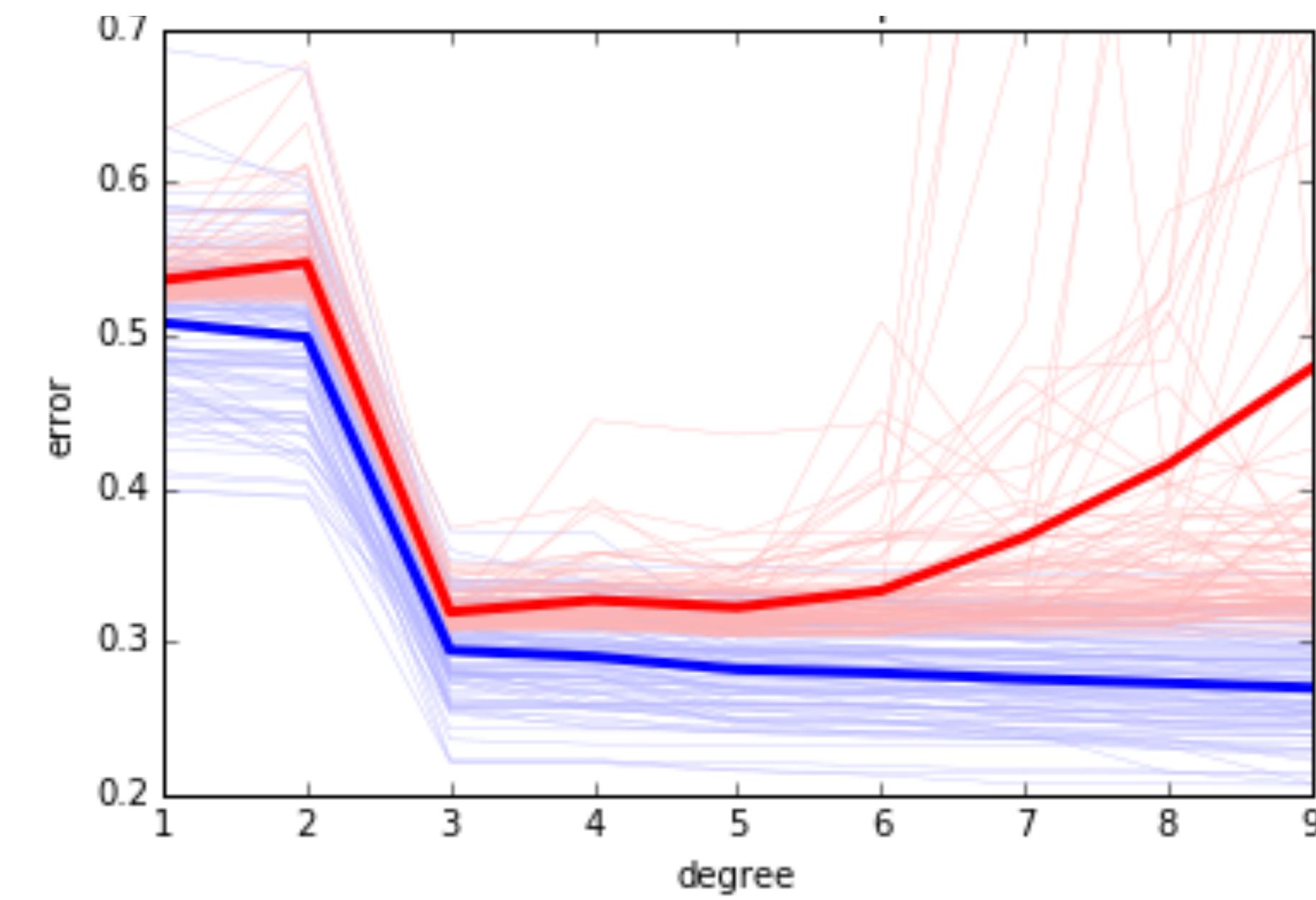
Draw learning curve!  
(esp. easy w/ SGD)



# Model selection curves



Ridge regression



Degree in case of a polynomial feature expansion

**But this depends on the  
algorithm!**

# Double descent curve

