# Bias-Variance Decomposition

Machine Learning Course - CS-433
Oct 2, 2024
Martin Jaggi & Nicolas Flammarion

**EPFL**

# Last time

How can we judge if a given predictor is good?
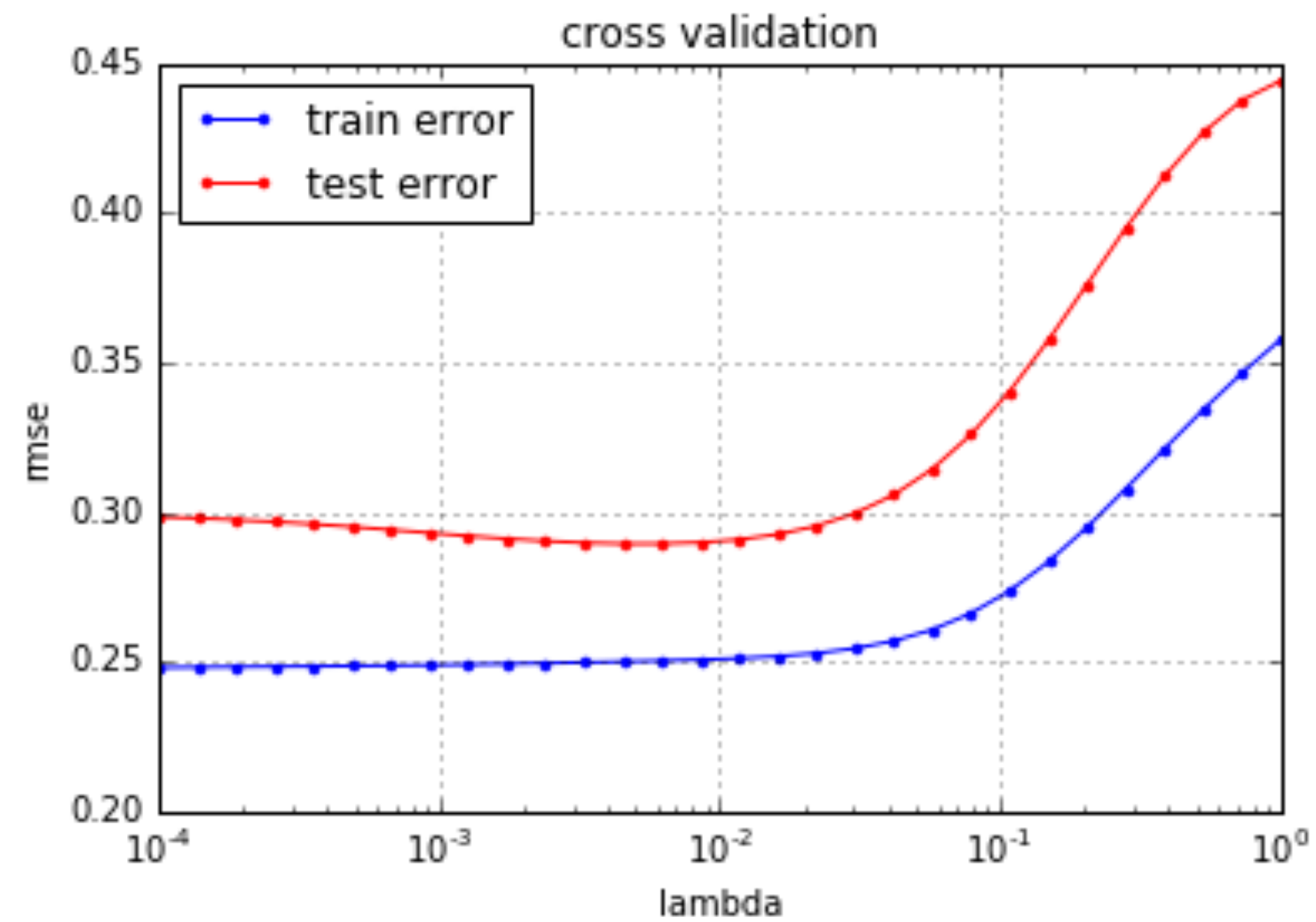How to select the best models of a family?

➡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)
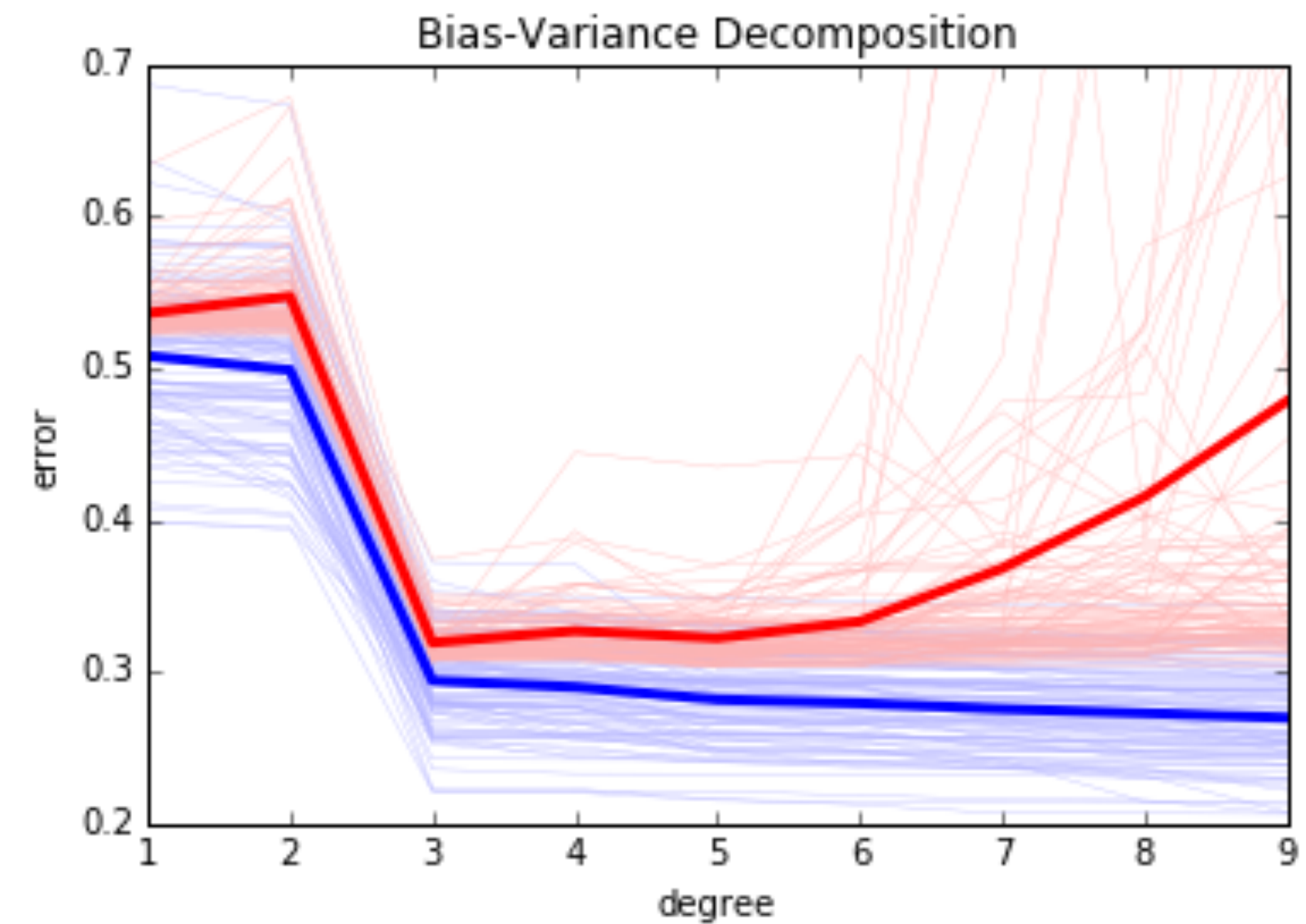
Motivation: Hyperparameters search (which often control the complexity)

But we haven't investigated the role of the complexity of the class
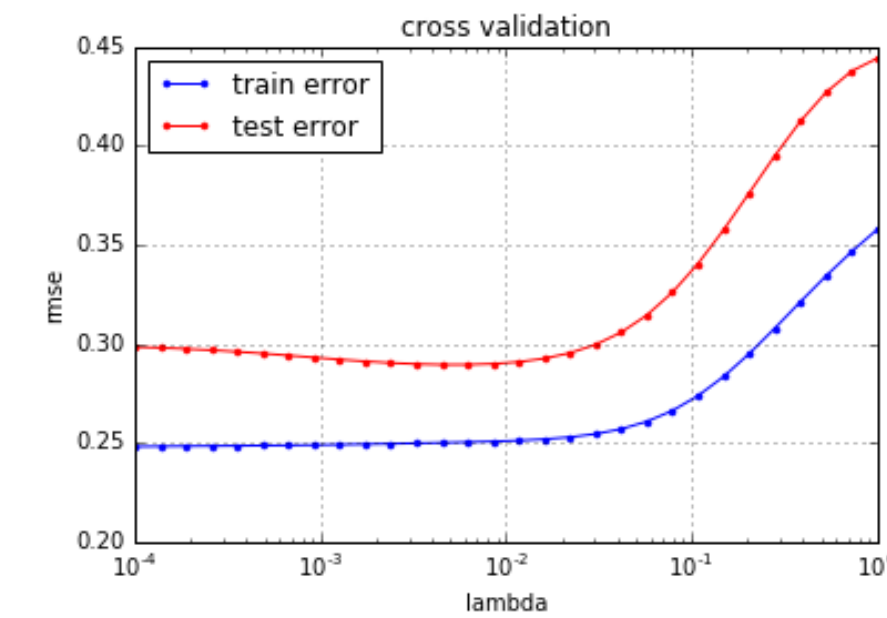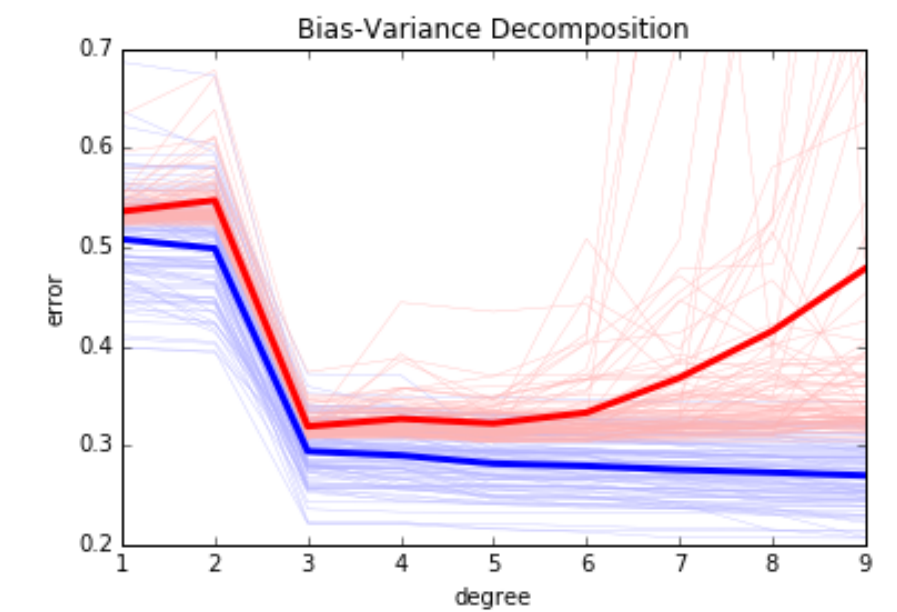
# Model selection curves



Ridge regression



Degree in case of a polynomial feature expansion

# Today



Ridge regression



Polynomial feature expansion

How does the risk behave as a function of the complexity of the model class?
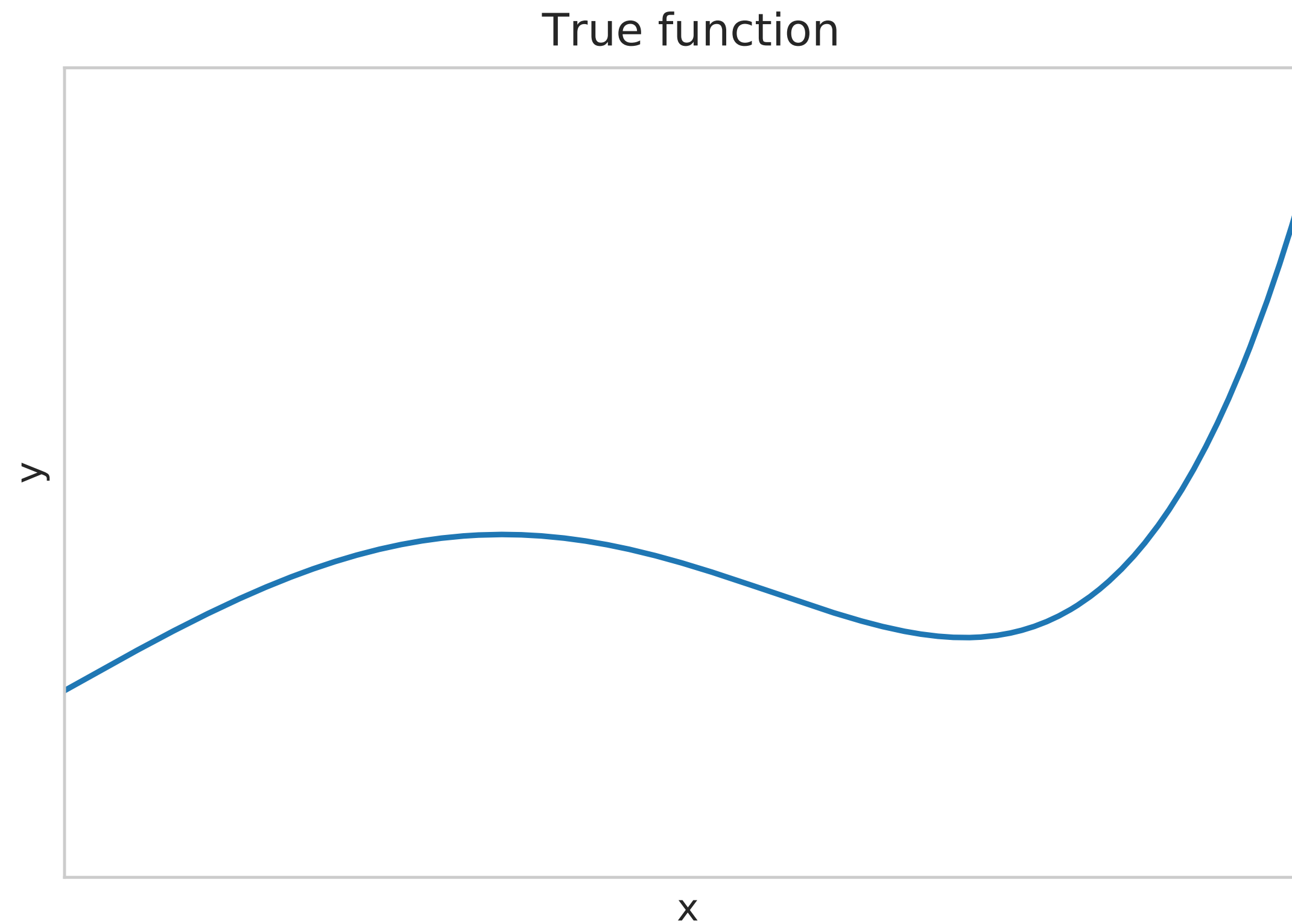
➡ ***Bias-Variance tradeoff***

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

⚠ Before: quantitative

Now: ***qualitative***

# A small experiment: 1D-regression

True function

# A small experiment: 1D-regression
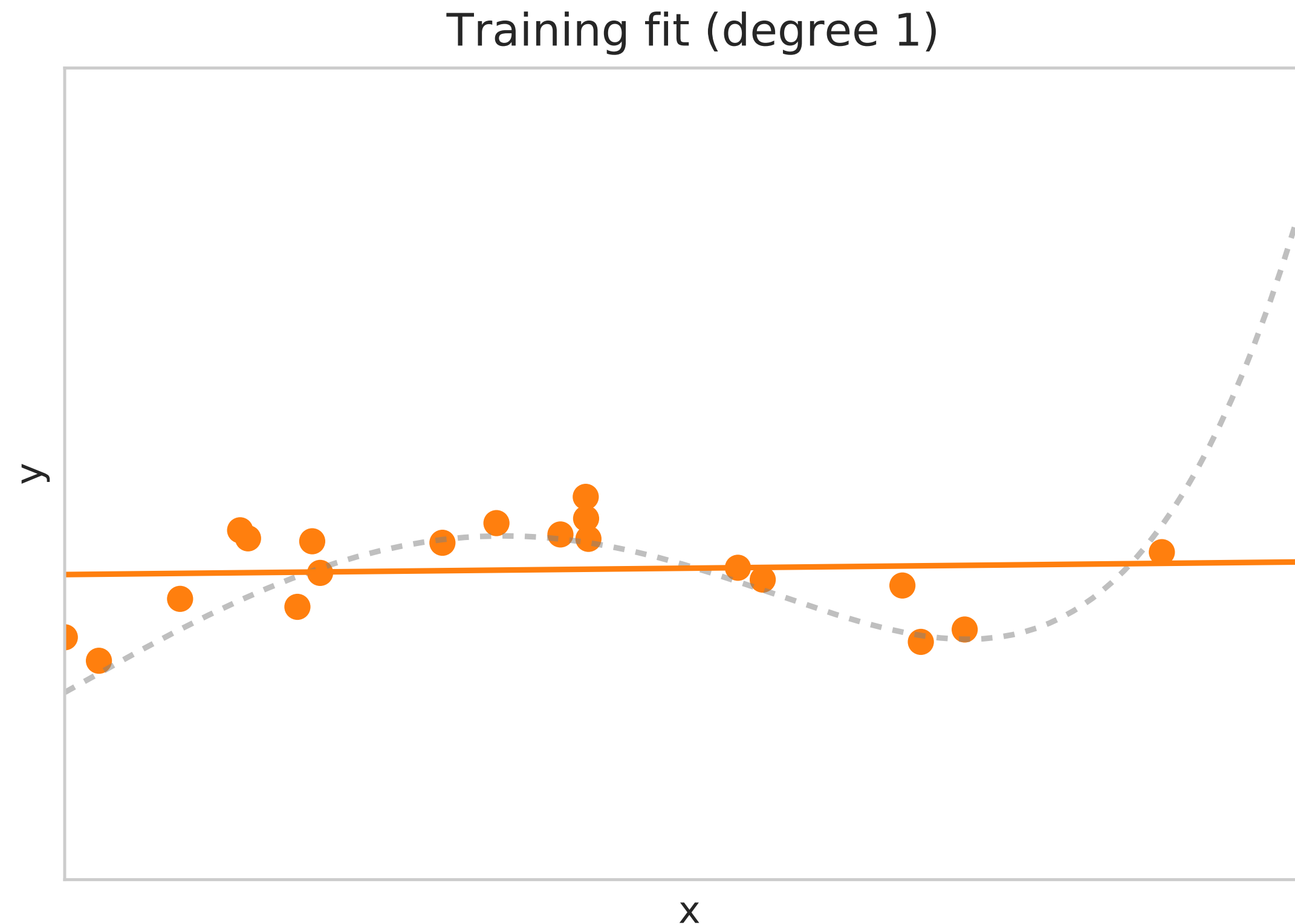
Sampled training set



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

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

➡ How far should you go?

# Simple model: bad fit



Training fit (degree 1)

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

# Complex model: good fit?



Training fit (degree 12)

High degree polynomial will be a good fit. But?

# But there is randomness in the data



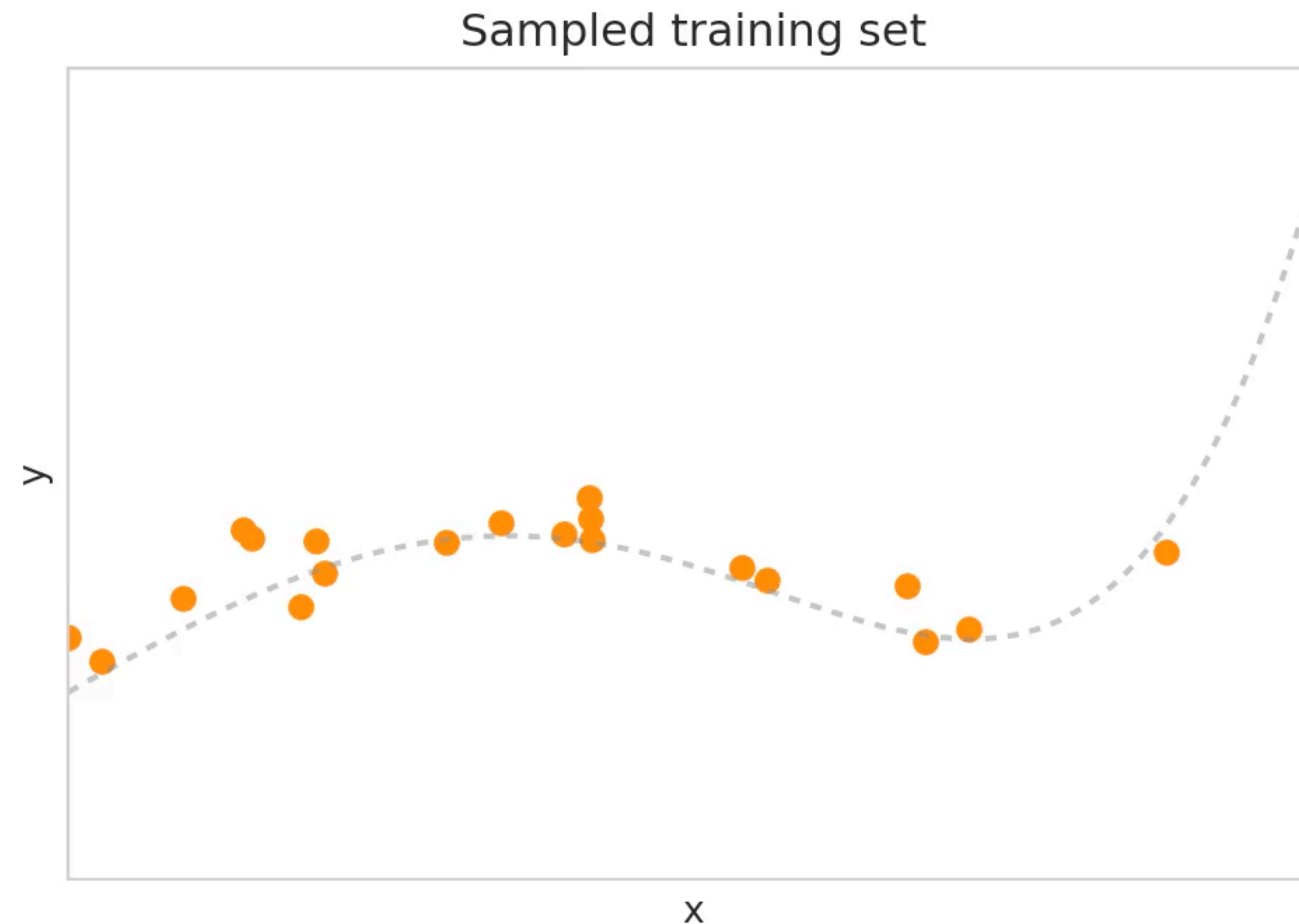Sampled training set

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

# But there is randomness in the data

Sampled training set



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

# Thus there is randomness in the predictions

Training fit (degree 1)

Training fit (degree 12)

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

Changing one of the observations may change the prediction considerably

**Underfitting**

**Overfitting**

## Simple models are less sensitive

# Simple models have large bias but low variance



Learned functions (degree 1)

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



Learned functions (degree 9)

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



Learned functions (degree 4)

# Data model: output perturbed by some noise

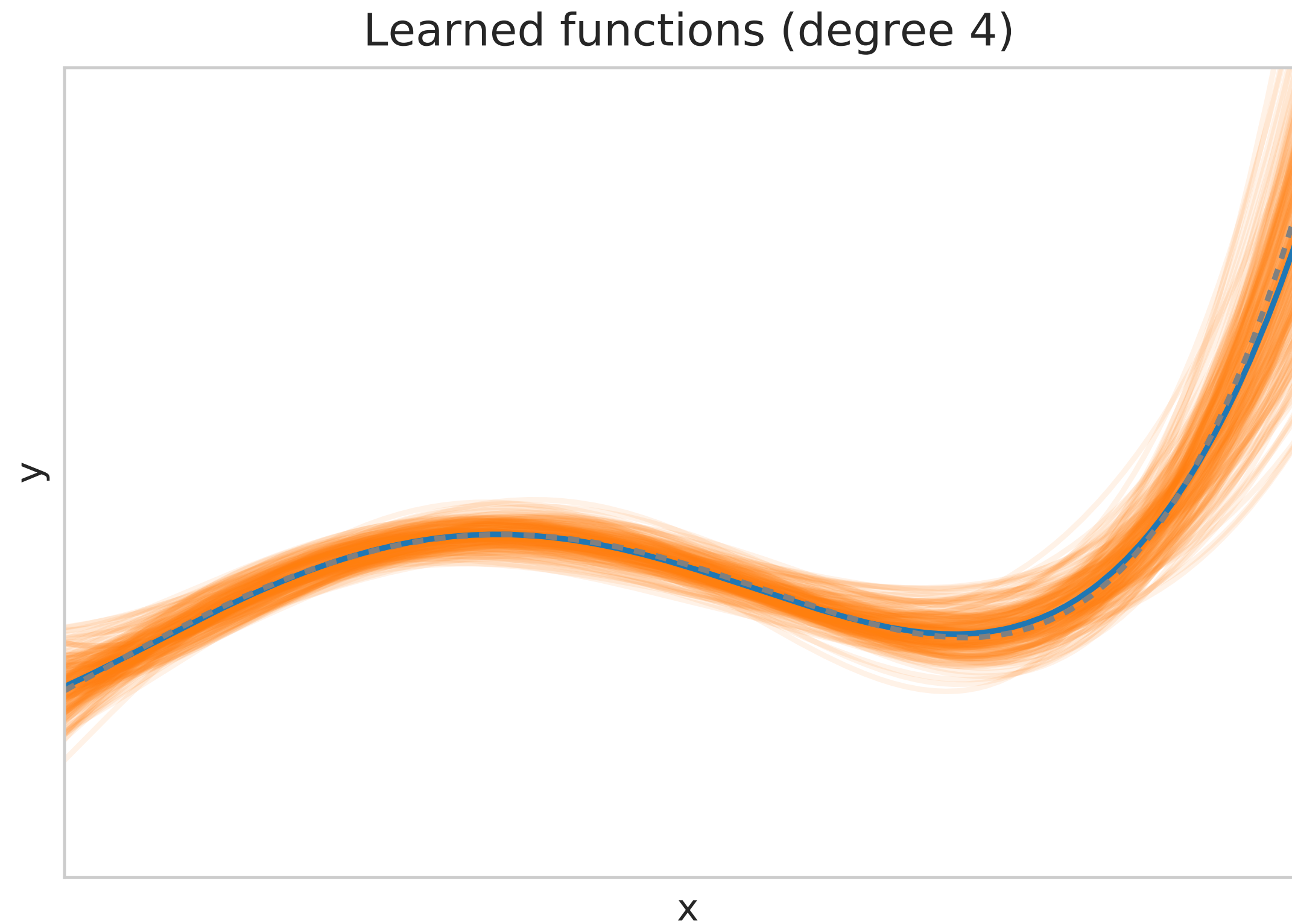**<u>Output</u>**

True model perturbed by some noise

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

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

**<u>True model</u>**

$f$ arbitrary and unknown function
The model is **not realizable**, i.e.,
$f$ is not in our model class

**<u>Input</u>**

$x \sim \mathscr{D}_x$

(fixed but unknown)

**<u>Noise</u>**

$\varepsilon \sim \mathscr{D}_\varepsilon$ 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

Train Data → Algorithm → Prediction

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

$$\mathscr{A}$$

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

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

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

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

# Error Decomposition



Train Data → Algorithm → Prediction
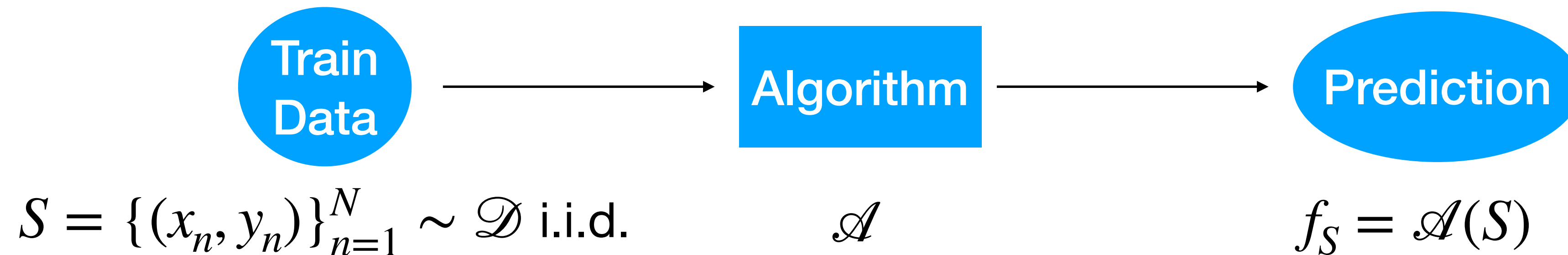
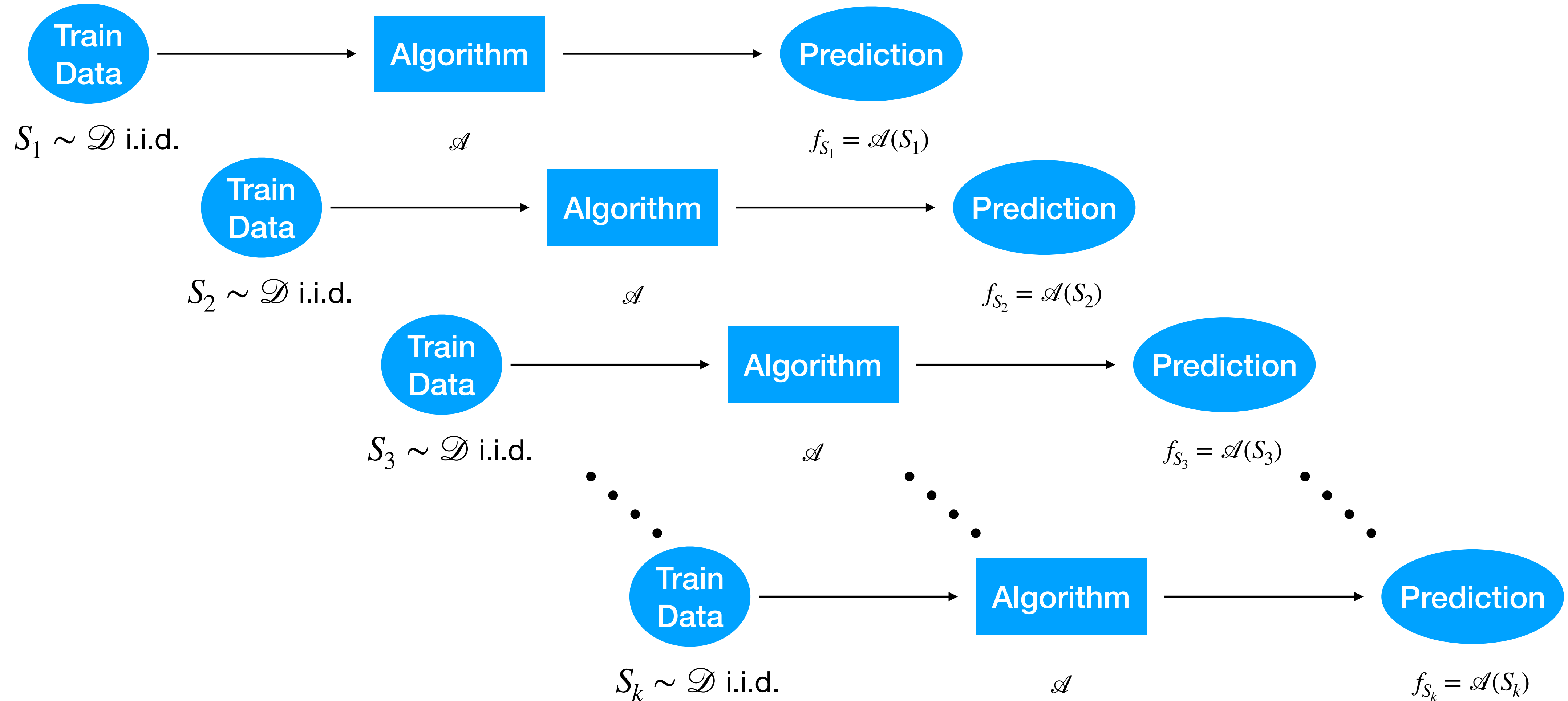$S = \{(x_n, y_n)\}_{n=1}^{N} \sim \mathcal{D}$ i.i.d.     $\mathcal{A}$     $f_S = \mathcal{A}(S)$

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 $x_0$:

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

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

# We run the experiment many times



$S_1 \sim \mathscr{D}$ i.i.d.    $\mathscr{A}$    $f_{S_1} = \mathscr{A}(S_1)$

$S_2 \sim \mathscr{D}$ i.i.d.    $\mathscr{A}$    $f_{S_2} = \mathscr{A}(S_2)$

$S_3 \sim \mathscr{D}$ i.i.d.    $\mathscr{A}$    $f_{S_3} = \mathscr{A}(S_3)$

$S_k \sim \mathscr{D}$ i.i.d.    $\mathscr{A}$    $f_{S_k} = \mathscr{A}(S_k)$

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

# A decomposition in three terms

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

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

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

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

First we expand the square:

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

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

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

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

- $\mathbb{E}_{\varepsilon\sim\mathscr{D}_\varepsilon}[\varepsilon^2] = \mathsf{Var}_{\varepsilon\sim\mathscr{D}_\varepsilon}[\varepsilon]$

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

Therefore

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

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 $S$

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

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

$$\left. + 2(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)])(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0))\right]$$

Cross-term:

$$\mathbb{E}_{S \sim \mathcal{D}}\left[\left(f(x_0) - \mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)]\right) \cdot \left(\mathbb{E}_{S' \sim \mathcal{D}}[f_{S'}(x_0)] - f_S(x_0)\right)\right]$$

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

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

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

# Bias-Variance Decomposition

We obtain the following decomposition into three positive terms:

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

$$\textbf{Bias} \quad \longrightarrow \quad + \quad (f(x_0) - \mathbb{E}_{S'\sim\mathscr{D}}[f_{S'}(x_0)])^2$$

$$\textbf{Variance} \quad \longrightarrow \quad + \quad \mathbb{E}_{S\sim\mathscr{D}}\big[(f_S(x_0) - \mathbb{E}_{S'\sim\mathscr{D}}[f_{S'}(x_0)])^2\big]$$

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

True Error ↑

Noise ┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄

→ Model Complexity

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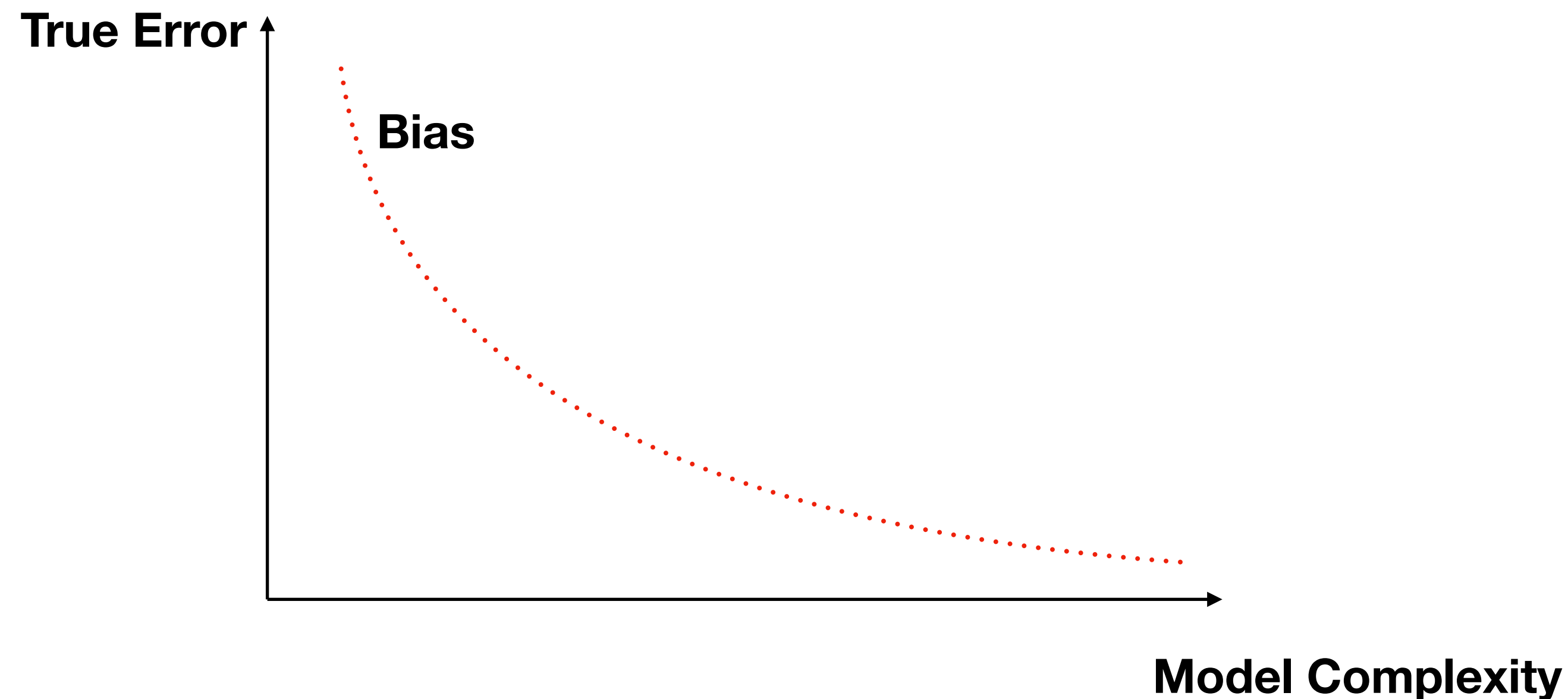- It is not possible to predict the noise from the data since they are independent

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



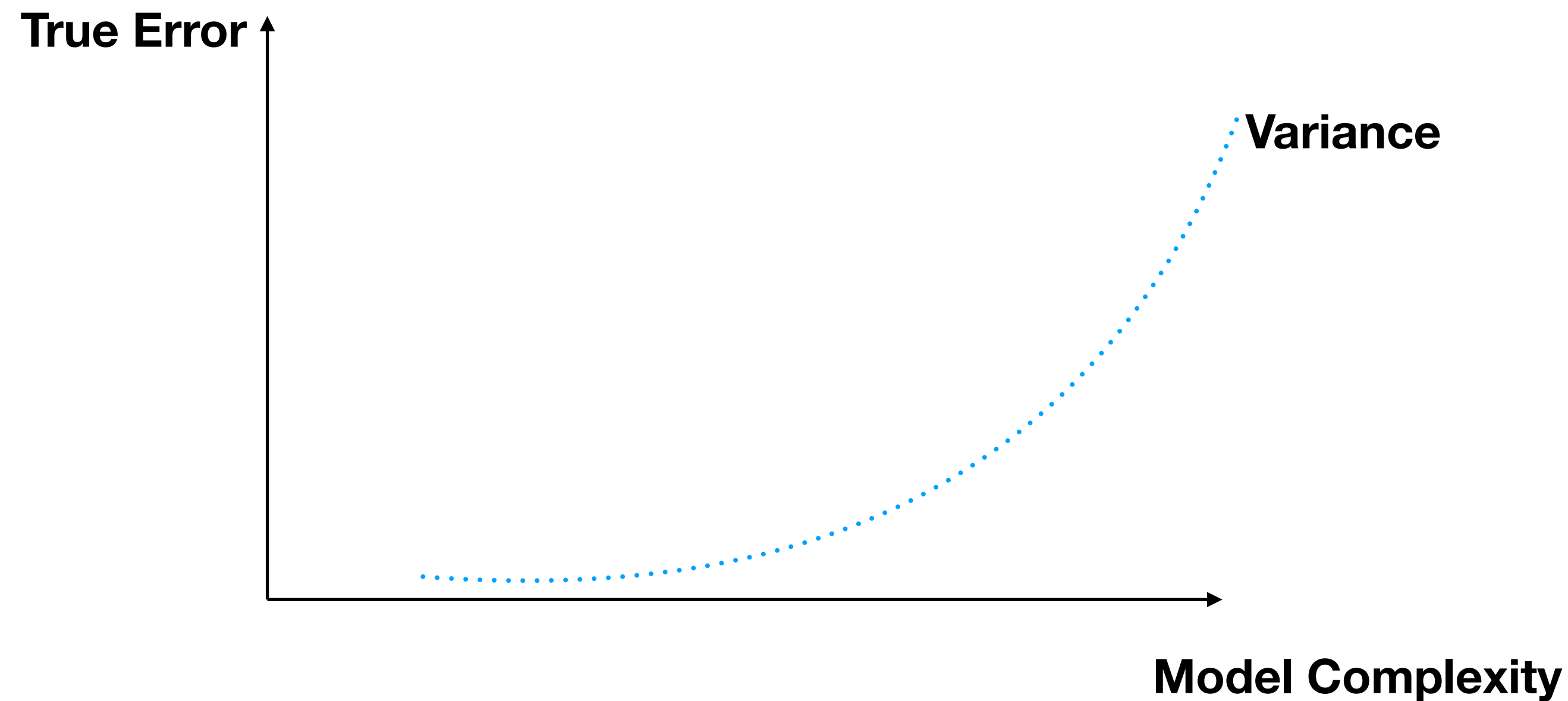Learned functions (degree 1)

Learned functions (degree 9)

- 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**

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



- 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**

# Variance: $\mathbb{E}_{S \sim \mathscr{D}}\left[(f_S(x_0) - \mathbb{E}_{S \sim \mathscr{D}}[f_S(x_0)])^2\right]$



Learned functions (degree 1)

Learned functions (degree 9)

- 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
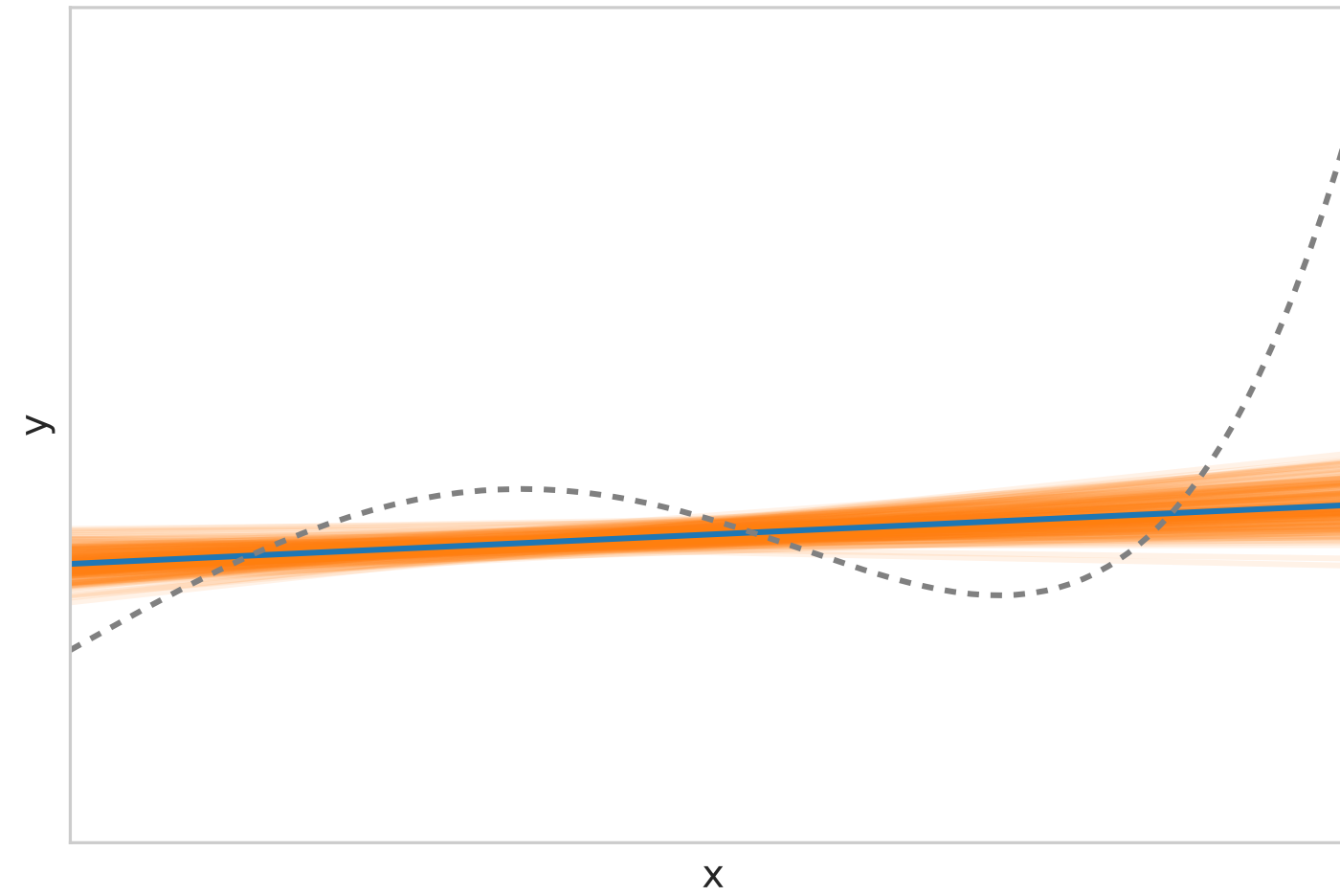
# Variance: $\mathbb{E}_{S\sim\mathscr{D}}\left[(f_S(x_0) - \mathbb{E}_{S\sim\mathscr{D}}[f_S(x_0)])^2\right]$



- 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



Learned functions (degree 1)
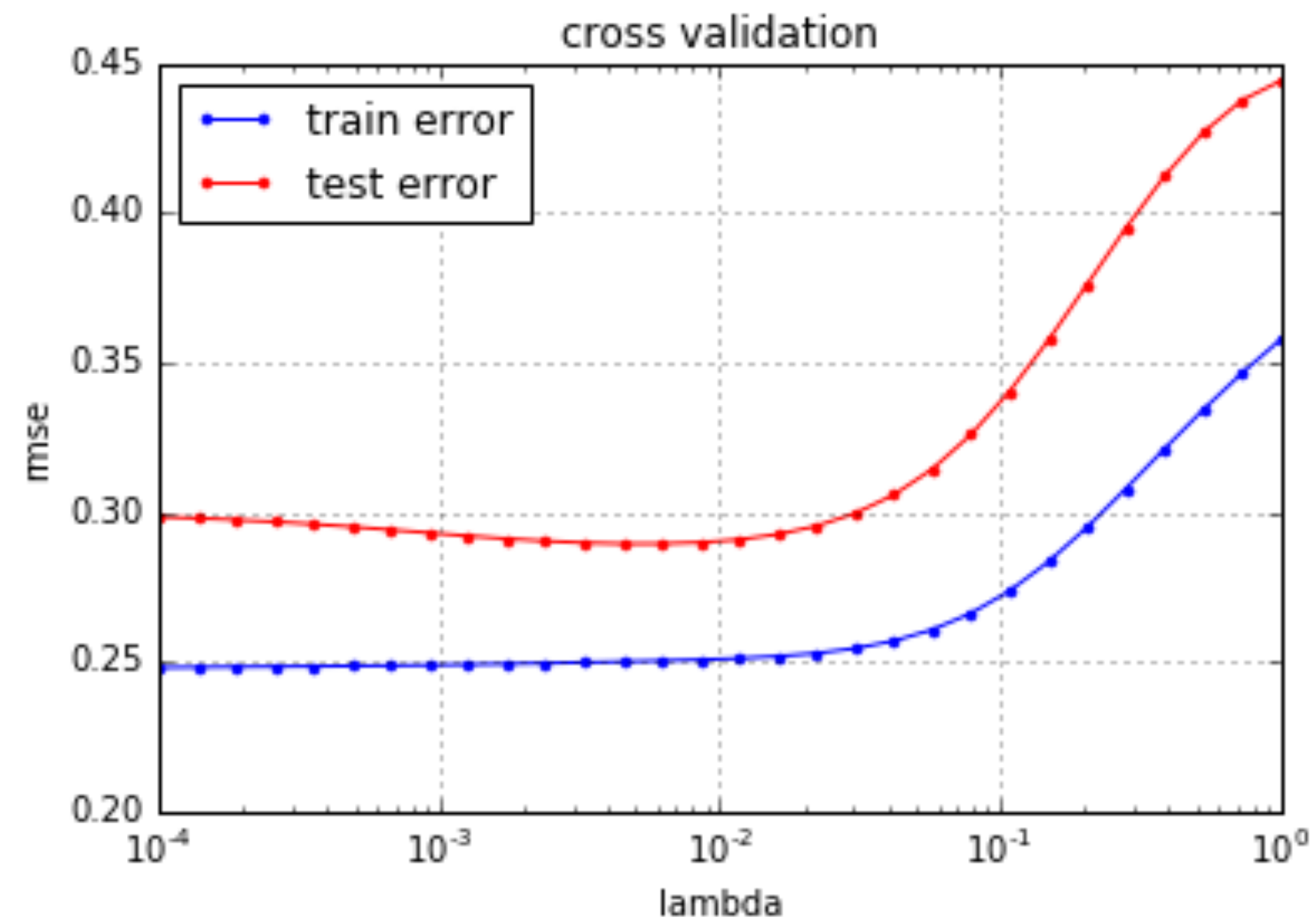
Learned functions (degree 4)
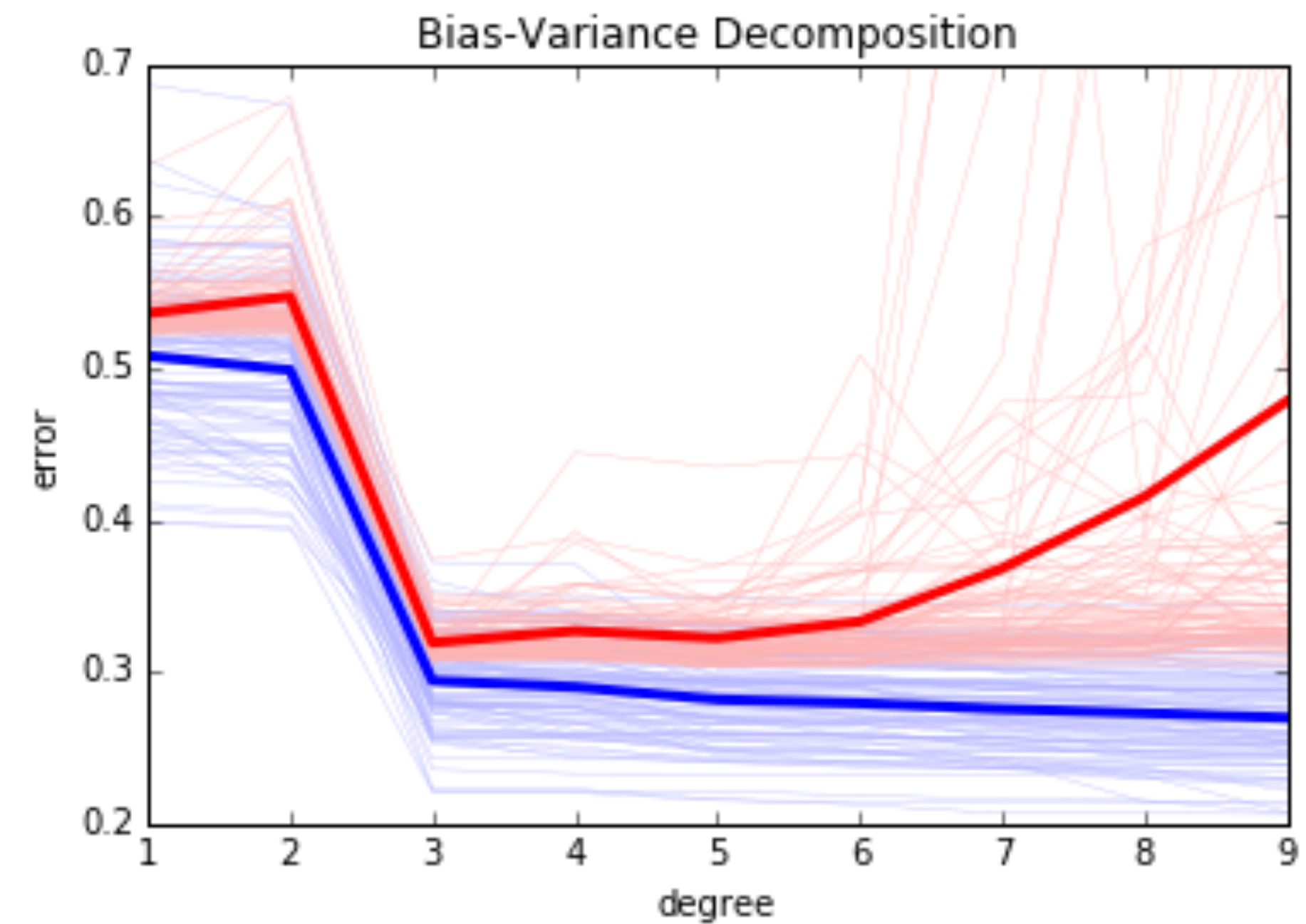
Learned functions (degree 9)
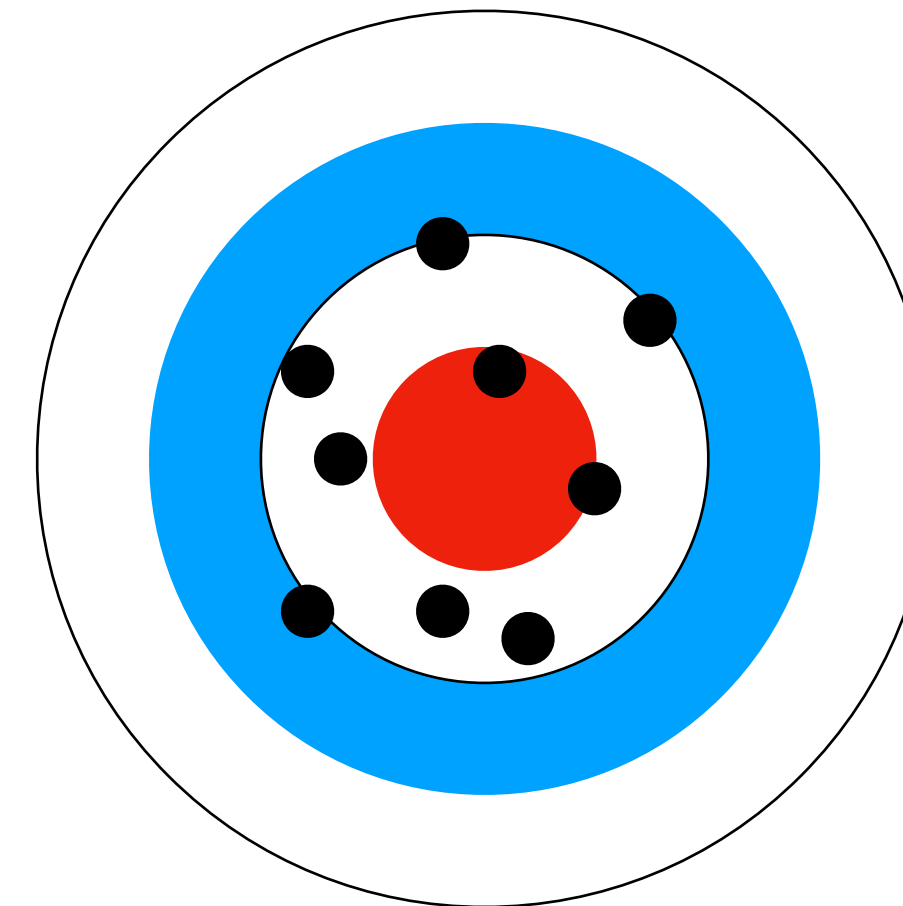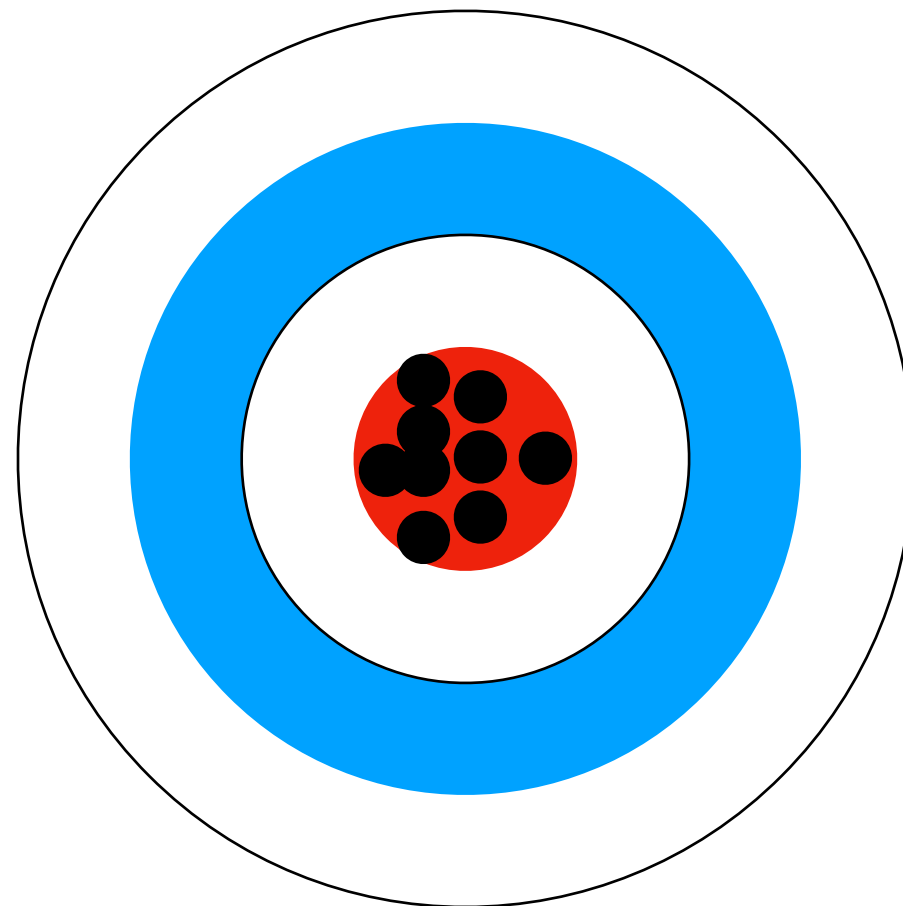
# Model selection curves



Ridge regression

Degree in case of a polynomial feature expansion

# Conclusion

# But this depends on the algorithm!
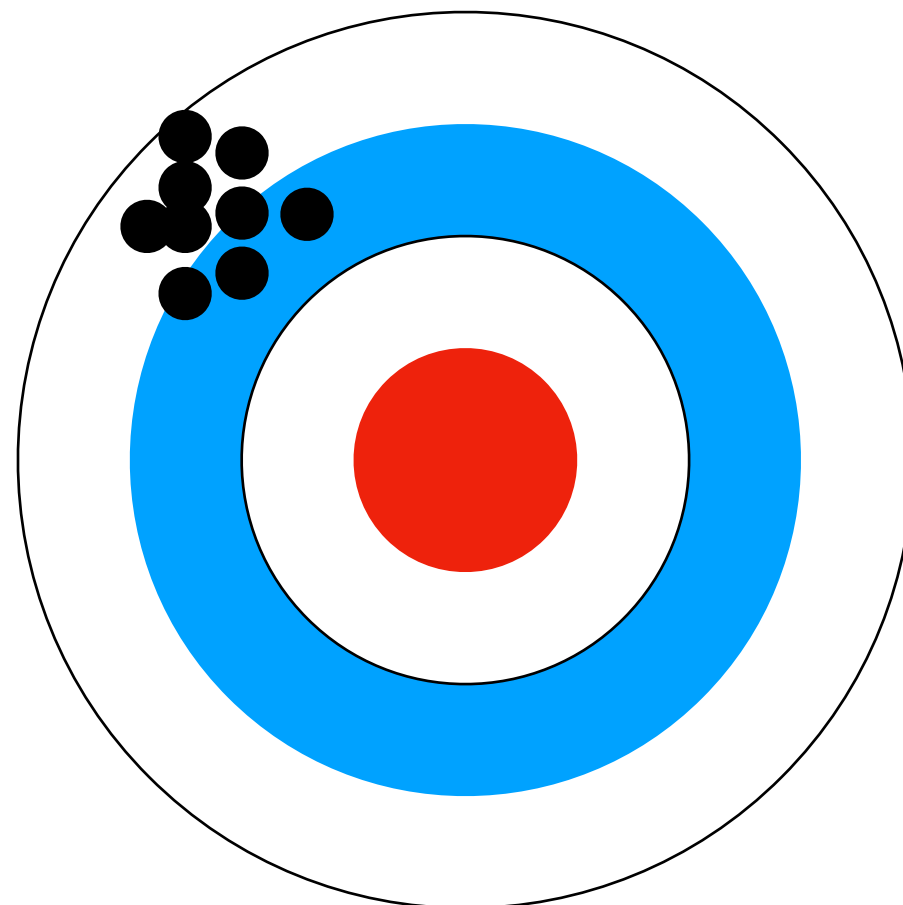
# Double descent curve