# Lecture 08

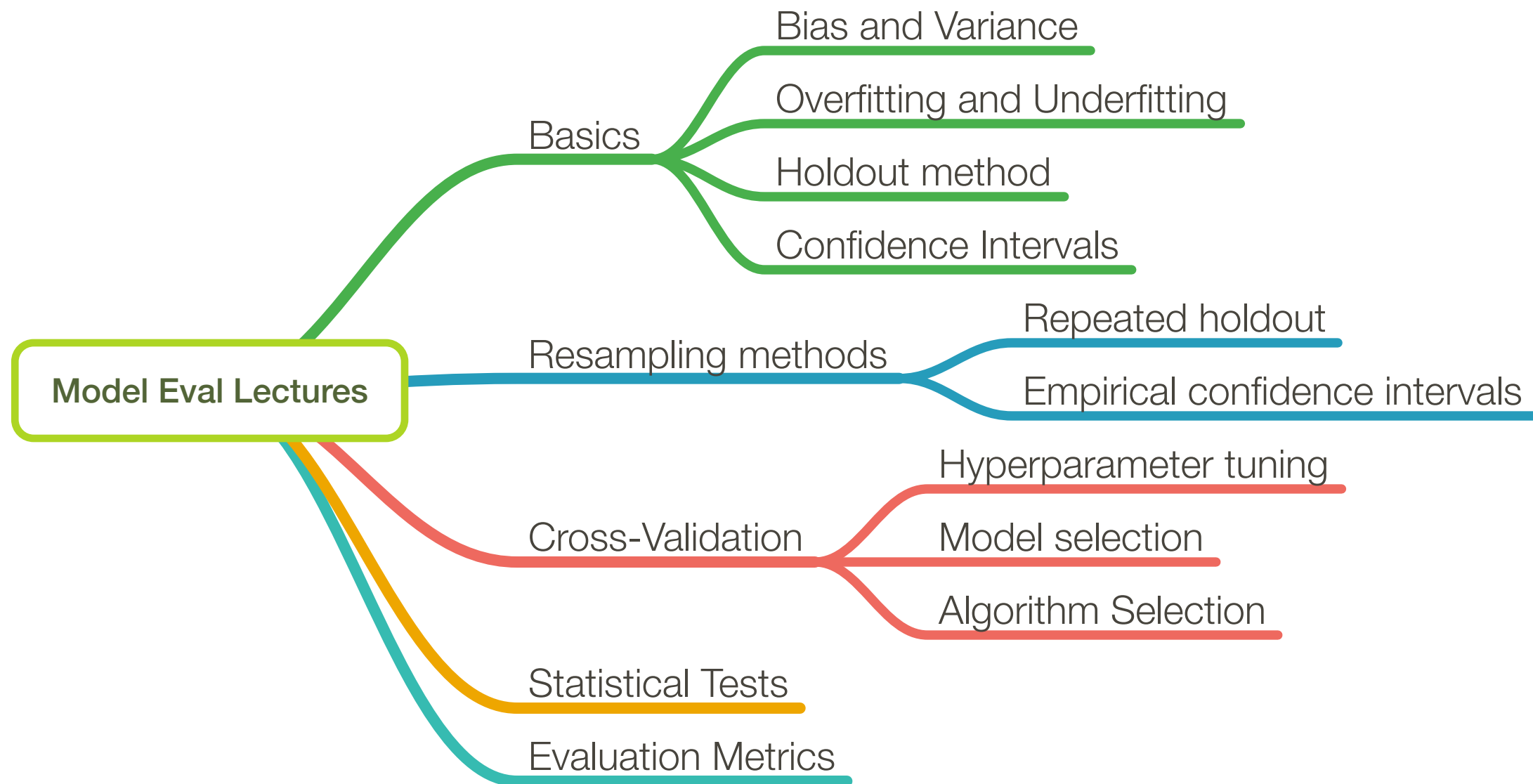# Model Evaluation 1:
## Introduction to Overfitting and Underfitting

STAT 479: Machine Learning, Fall 2018

Sebastian Raschka

http://stat.wisc.edu/~sraschka/teaching/stat479-fs2018/

# Overview

# Overfitting and Underfitting

# Overfitting and Underfitting

## "Generalization Performance"

- Goal is to fit a model that performs well on unseen inputs, that is, a model that "generalizes" well to unseen data

- A model that performs well on unseen inputs has a good generalization performance

- We say that a model with a good generalization performance has a "high generalization accuracy" or "low generalization error"
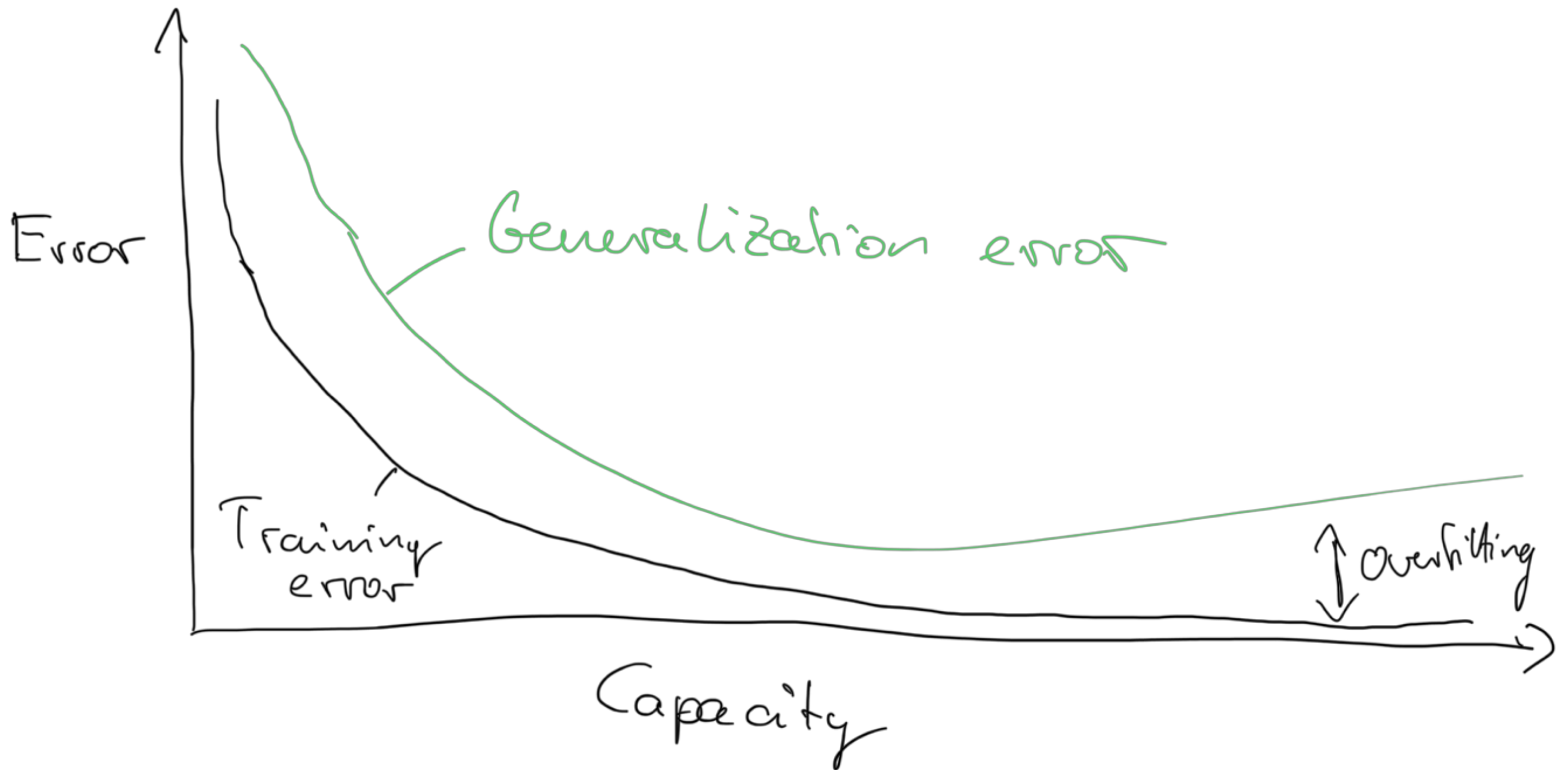
# Overfitting and Underfitting

## Assumptions

- i.i.d. assumption: inputs are independent, and training and test examples are identically distributed (drawn from the same probability distribution)

- The training error or accuracy provides an (optimistically) biased estimate of the generalization performance

- For some random model that has not been fitted to the training set, we expect both the training and test error to be equal

# Overfitting and Underfitting

## Model Capacity

- Underfitting: both training and test error are large

- Overfitting: gap between training and test error (where test error is higher)

- Generally, the larger the hypothesis space being searched by a learning algorithm, the higher its tendency to overfit (the size of the hypothesis space is related to the so-called "capacity" of a model); vice versa, models with small capacity do not even fit the training set well
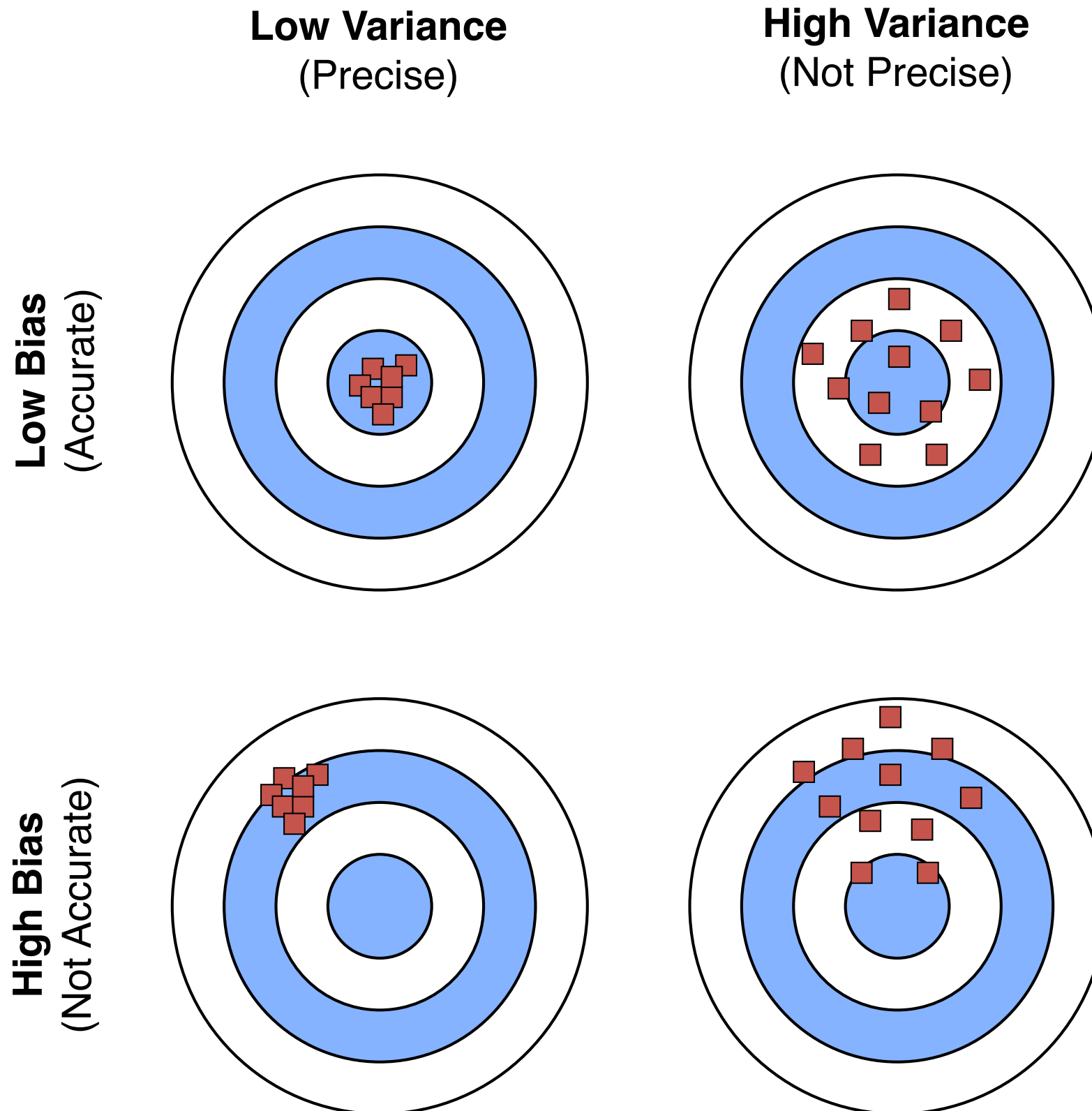
# Overfitting and Underfitting
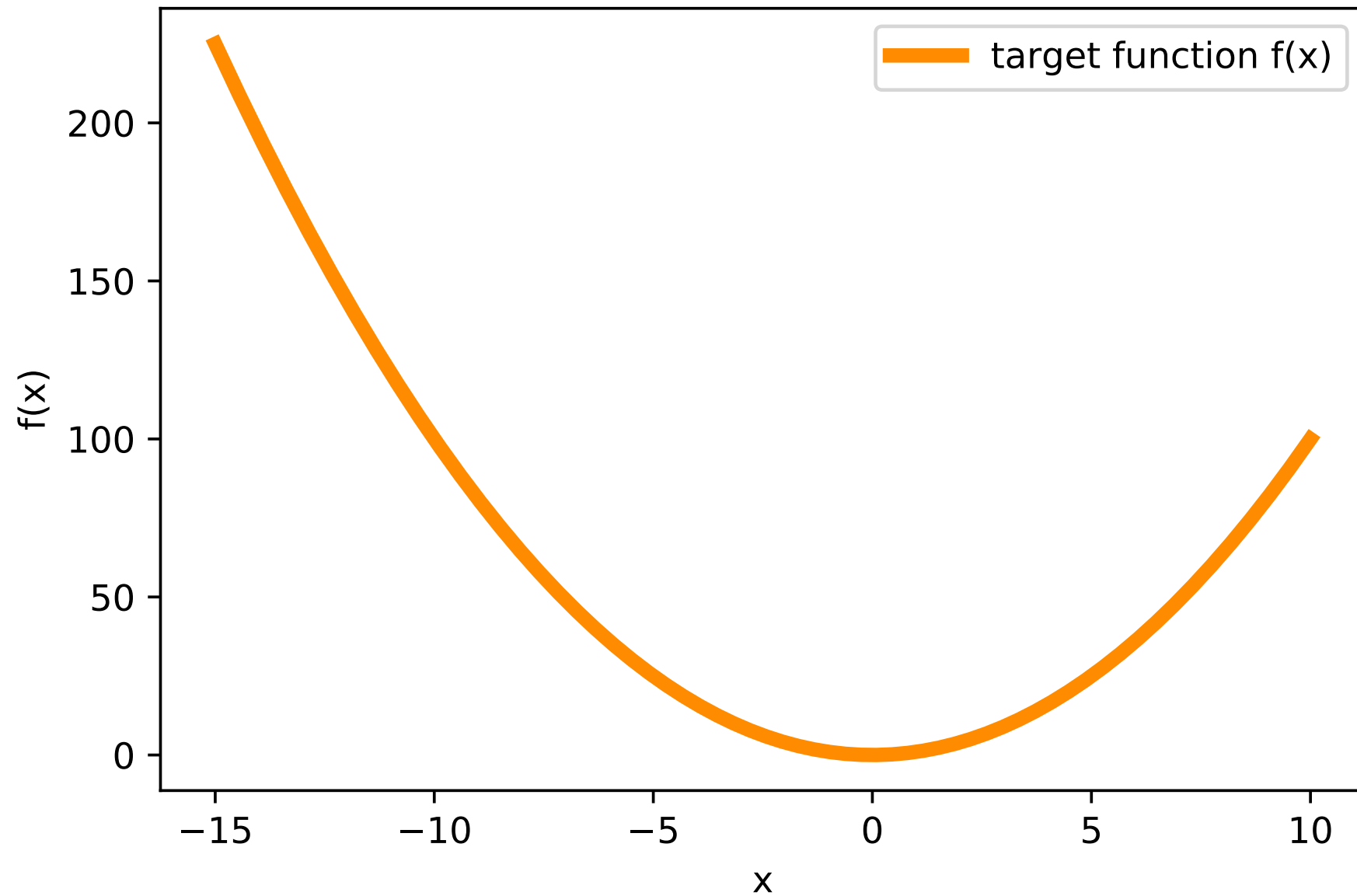
# Bias-Variance Decomposition and Trade-off

# Bias-Variance Decomposition

- Decomposition of the loss into bias and variance help us understand learning algorithms, concepts are correlated to underfitting and overfitting
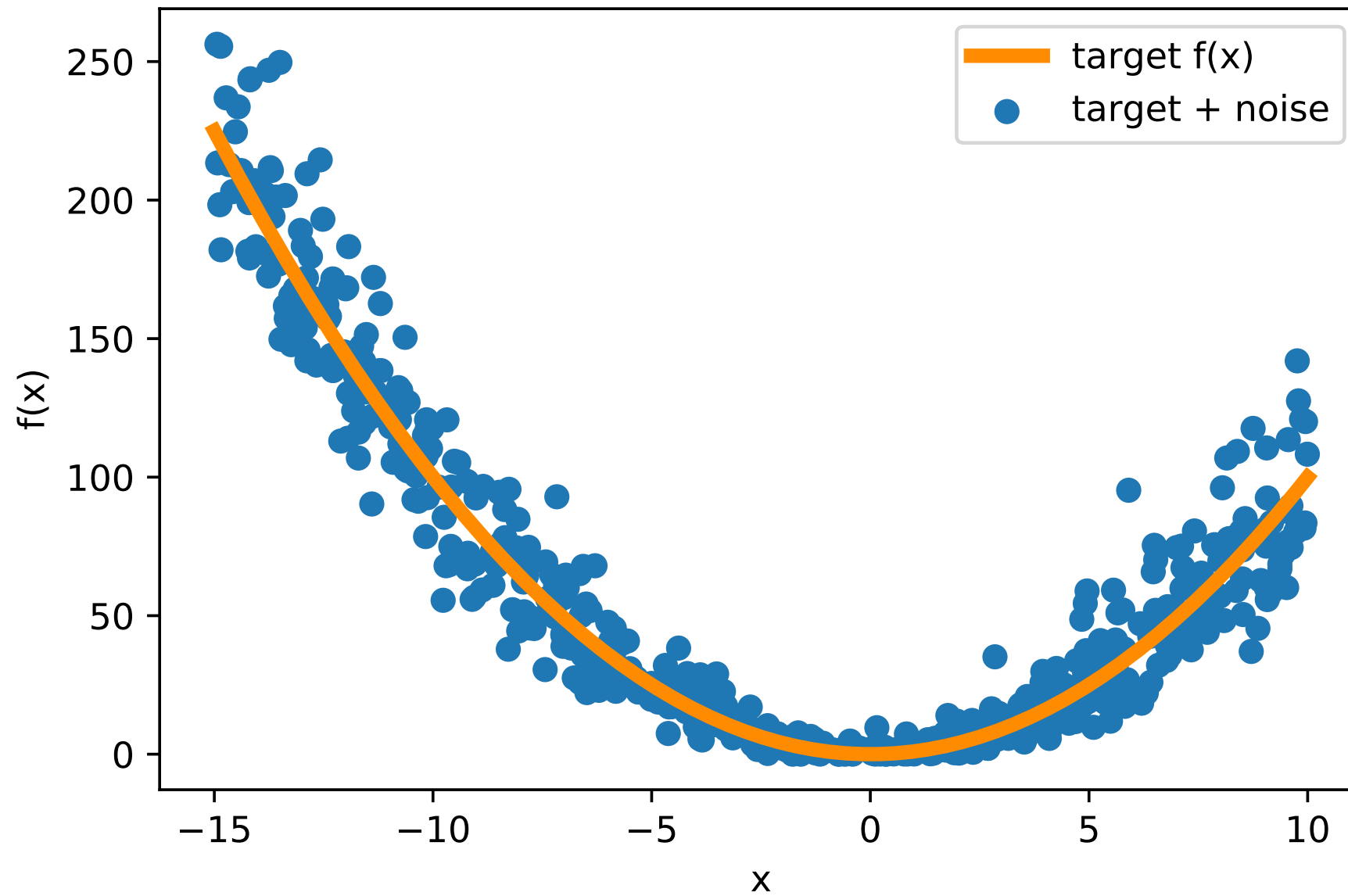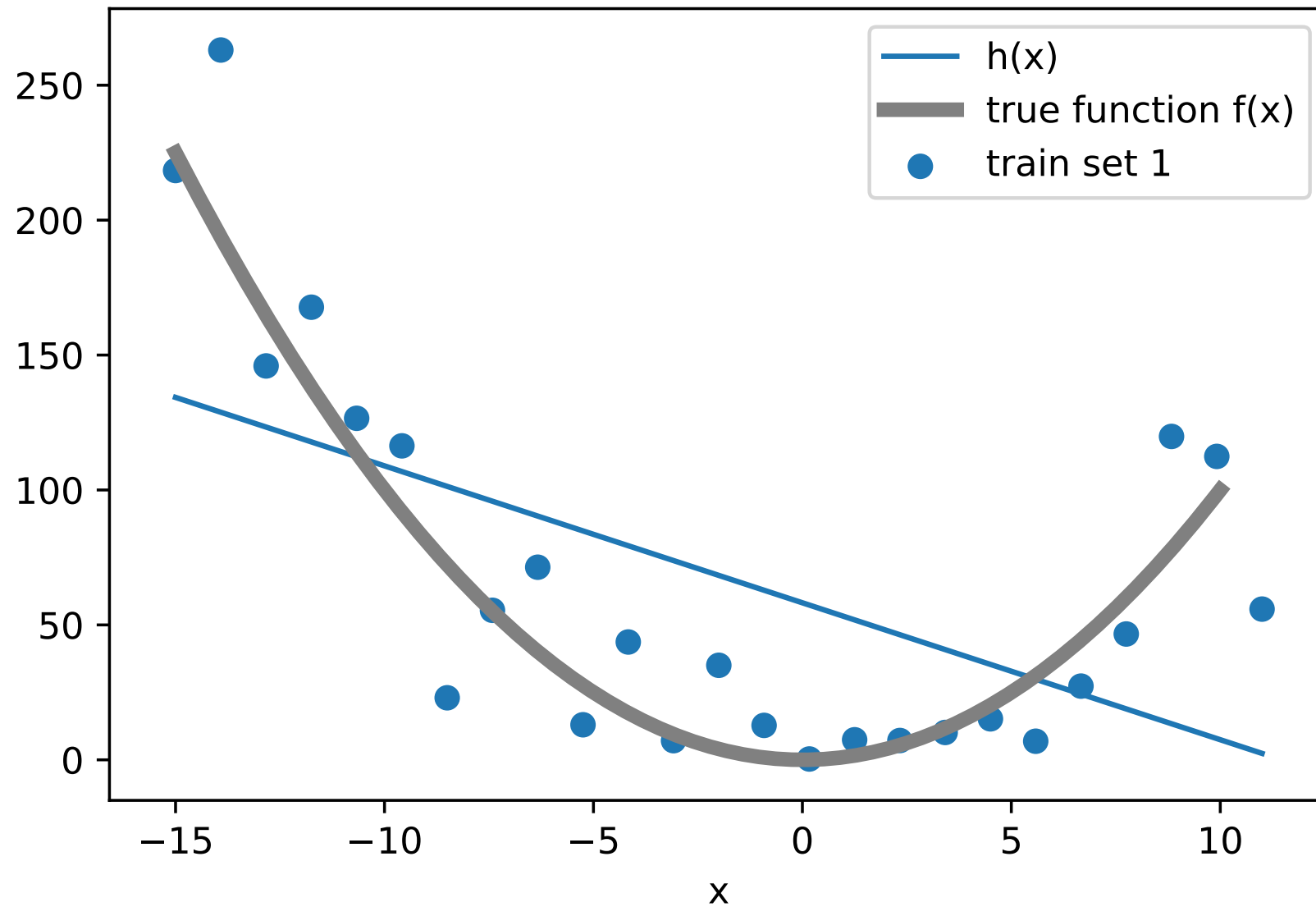
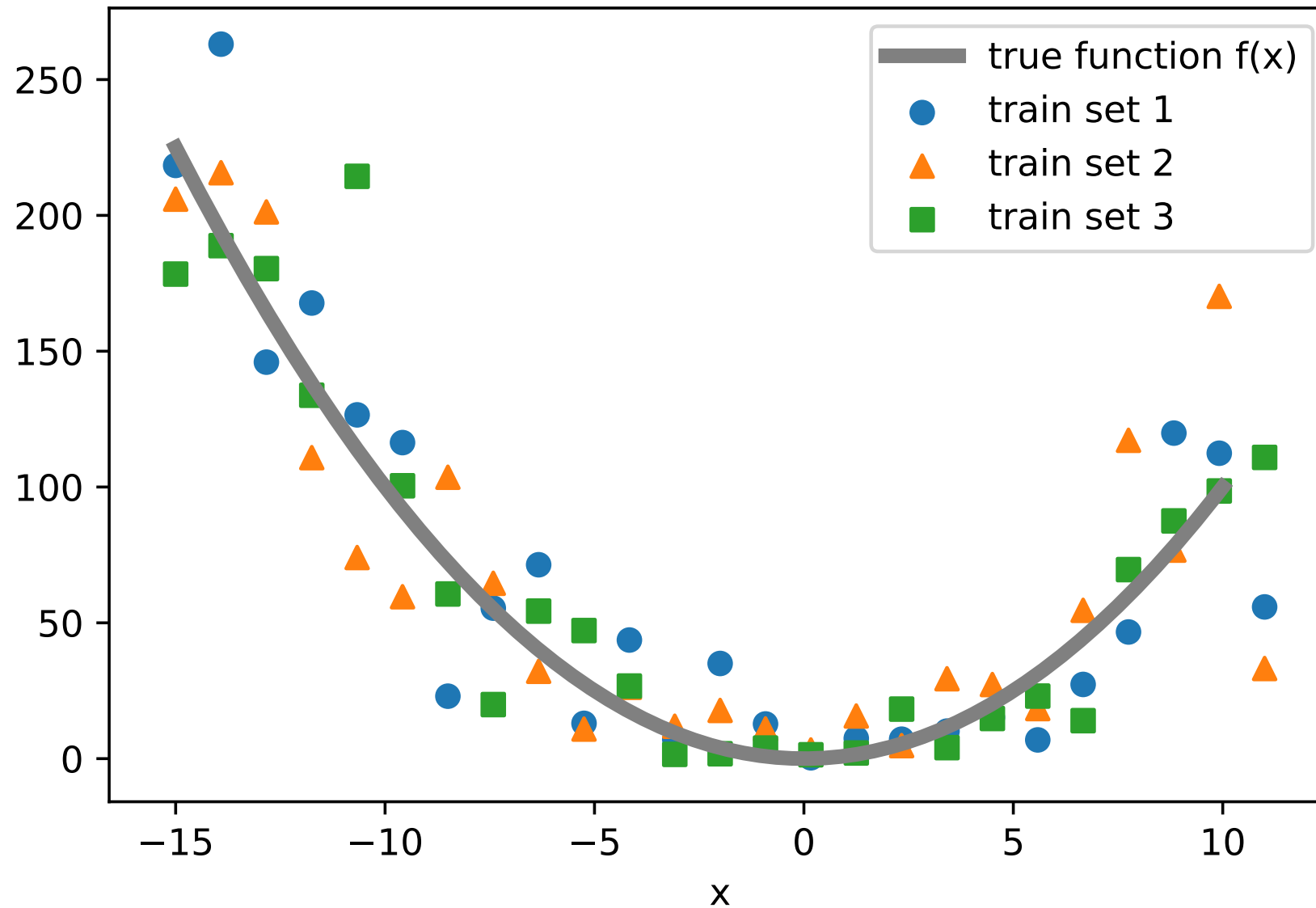# Bias-Variance Intuition

# Bias and Variance Intuition
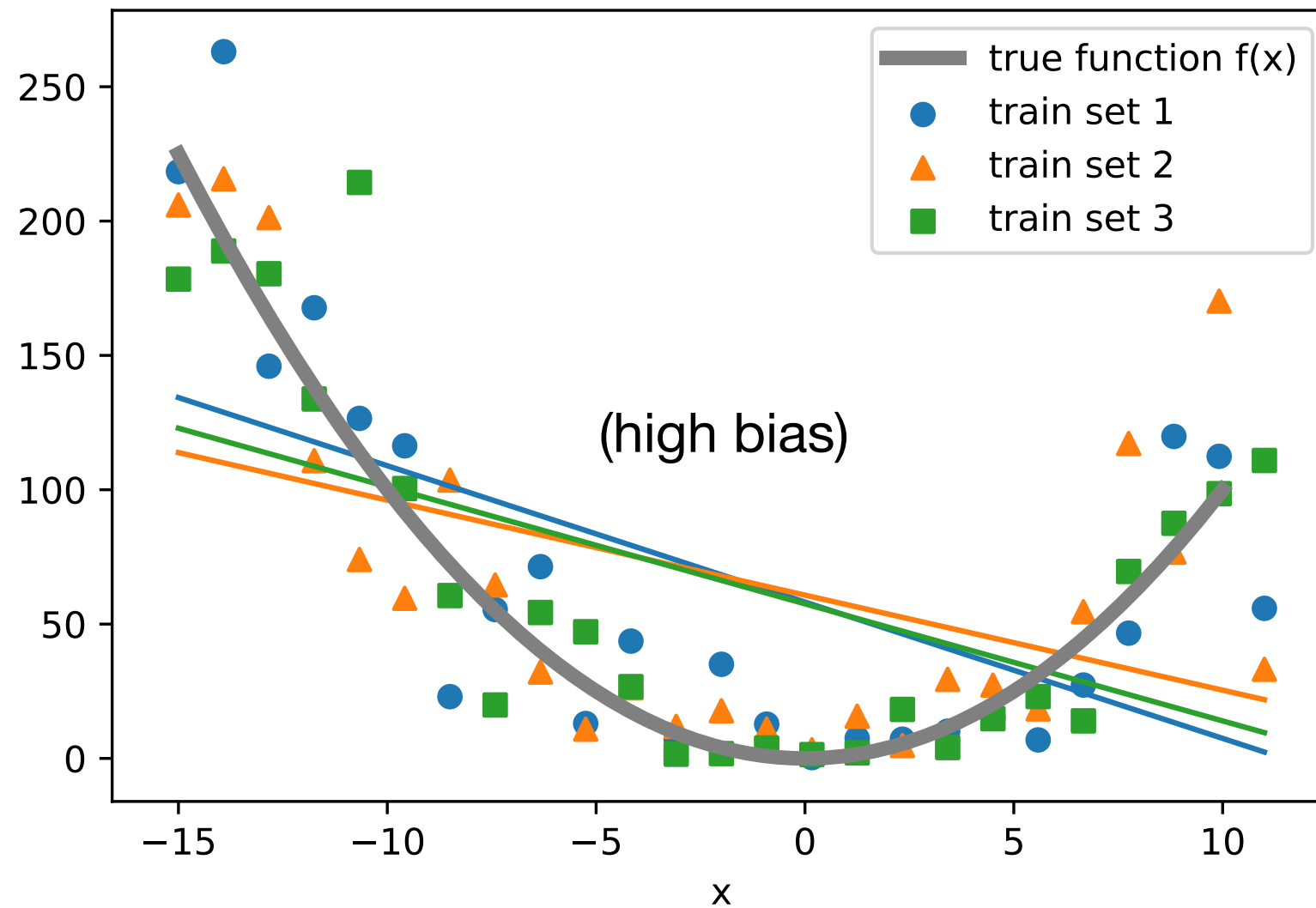
# Bias and Variance Intuition
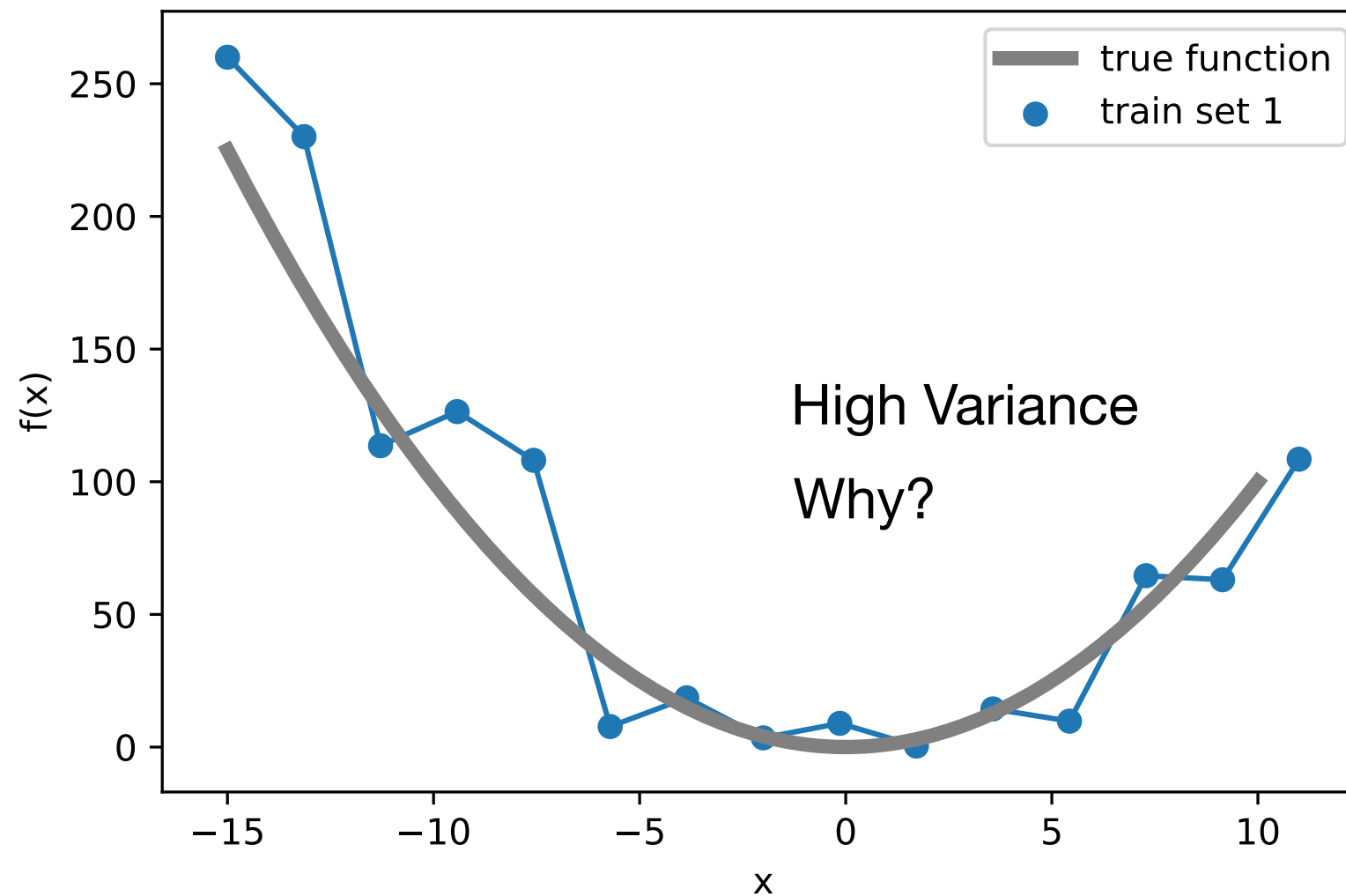
# Bias and Variance Intuition

# Bias and Variance Intuition

# Bias and Variance Intuition

# Bias and Variance Intuition



(here, I fit an unpruned decision tree)

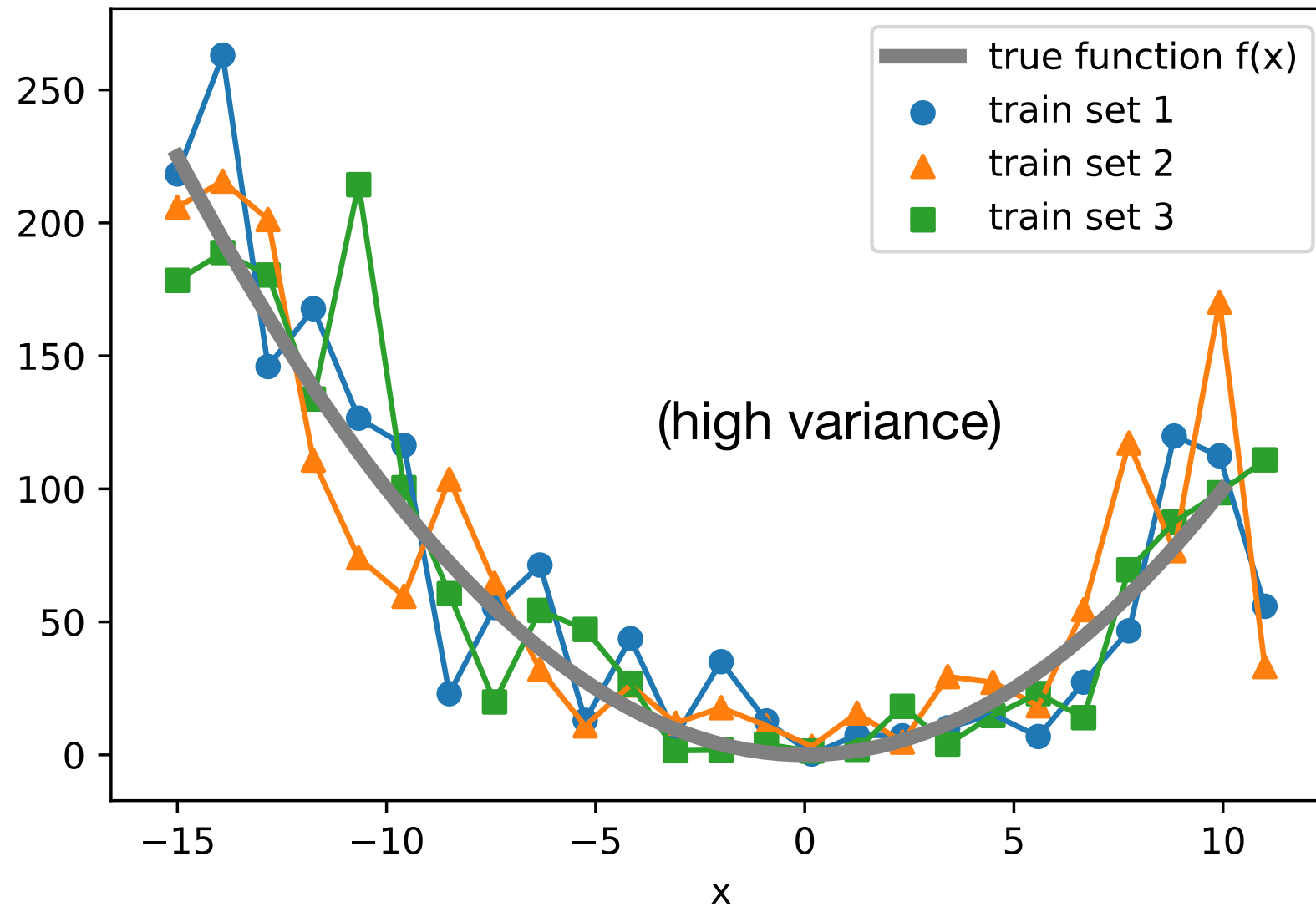# Bias and Variance Example



where f(x) is some true (target) function

suppose we have multiple training sets

# Bias and Variance Example



(high variance)

# Bias-Variance Decomposition

Point estimator:  $\hat{\theta} = f(x^{[1]}, x^{[2]}, \ldots, x^{[n]})$

of some parameter  $\theta$

(could also be a function, e.g., the hypothesis is an estimator of some target function)

# Bias-Variance Decomposition

Point estimator: $\quad \hat{\theta} = f(x^{[1]}, x^{[2]}, \ldots, x^{[n]})$

of some parameter $\theta$

(could also be a function, e.g., the hypothesis is
an estimator of some target function)

**Bias** $= E[\hat{\theta}] - \theta$

(the expectation is over the training data, i.e, the average estimator from
different training samples)

# Bias-Variance Decomposition of Squared Error

## General Definition:

### Intuition:

$$\mathbf{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

Bias is the difference between the average estimator from different training samples and the true value. (The expectation is over the training sets.)

$$\mathbf{Var}(\hat{\theta}) = E[\hat{\theta}^2] - \left( E[\hat{\theta}] \right)^2$$

The variance provides an estimate of how much the estimate varies as we vary the training data (e.g,. by resampling).

$$\mathbf{Var}(\hat{\theta}) = E[(E[\hat{\theta}] - \hat{\theta})^2]$$
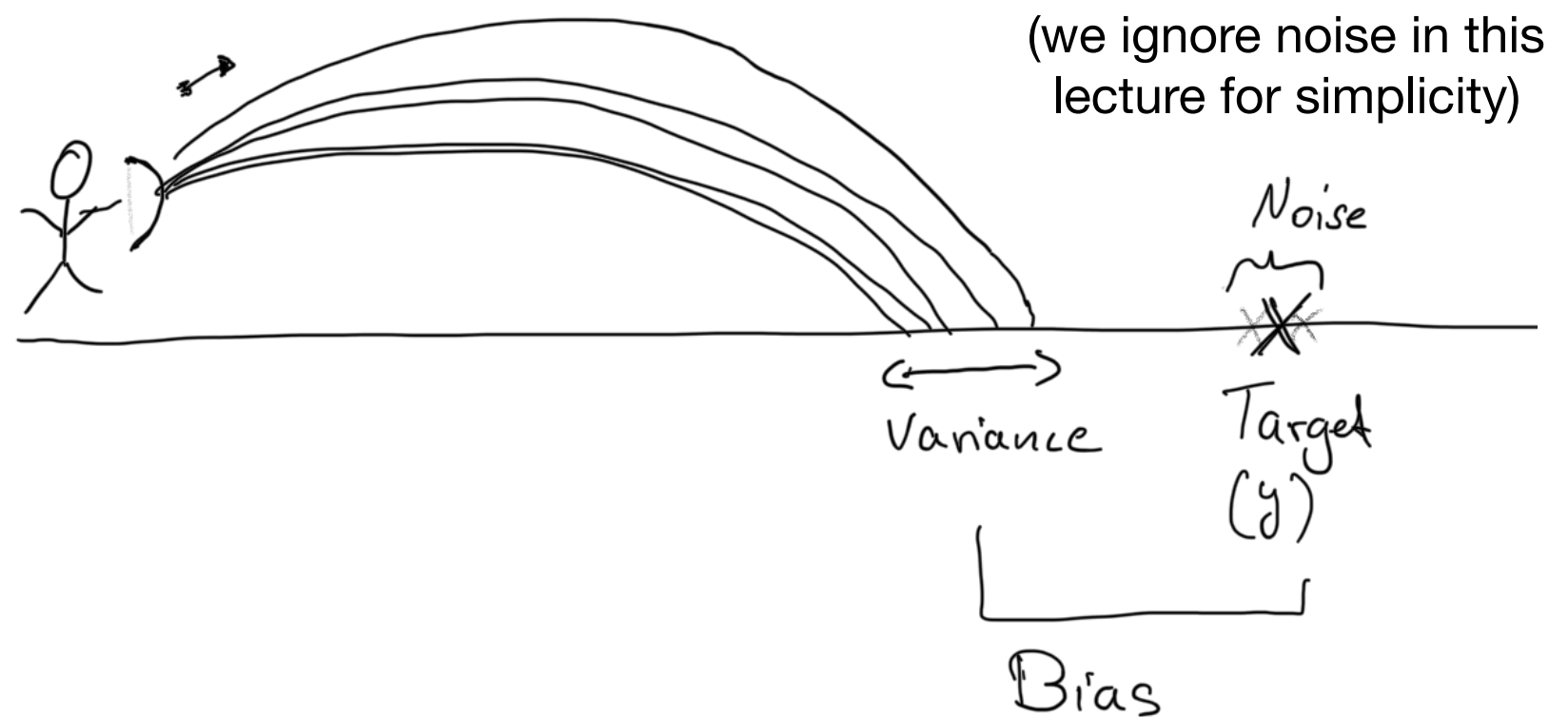
# Bias-Variance Decomposition of Squared Error

## General Definition:

$$\textbf{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$\textbf{Var}(\hat{\theta}) = E[\hat{\theta}^2] - \left( E[\hat{\theta}] \right)^2$$

$$\textbf{Var}(\hat{\theta}) = E[(E[\hat{\theta}] - \hat{\theta})^2]$$

## Intuition:

(we ignore noise in this lecture for simplicity)

# Bias-Variance Decomposition of Squared Error

General Definition:

"ML notation" for the Squared Error Loss:

$$\textbf{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$y = f(x) \qquad \text{(target, target function)}$$

$$\hat{y} = \hat{f}(x) = h(x)$$

$$\textbf{Var}(\hat{\theta}) = E[\hat{\theta}^2] - \left( E[\hat{\theta}] \right)^2$$

$$S = (y - \hat{y})^2$$

$$\textbf{Var}(\hat{\theta}) = E[(E[\hat{\theta}] - \hat{\theta})^2]$$

# Bias-Variance Decomposition of Squared Error

"ML notation" for the Squared Error Loss:

$$y = f(x) \quad \text{(target, target function)}$$

$$\hat{y} = \hat{f}(x) = h(x)$$

$$S = (y - \hat{y})^2$$

---

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - y)^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

# Bias-Variance Decomposition of Squared Error

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - y)^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

$$E[S] = E[(y - \hat{y})^2]$$

$$E[(y - \hat{y})^2] = (y - E[\hat{y}])^2 + E[(E[\hat{y}] - \hat{y})^2]$$

$$= \textbf{[Bias of the fit]}^2 + \textbf{Variance of the fit}$$

# Bias-Variance Decomposition of Squared Error

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - y)^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

???

$$E[S] = E[(y - \hat{y})^2]$$

$$E[(y - \hat{y})^2] = (y - E[\hat{y}])^2 + E[(E[\hat{y}] - \hat{y})^2]$$

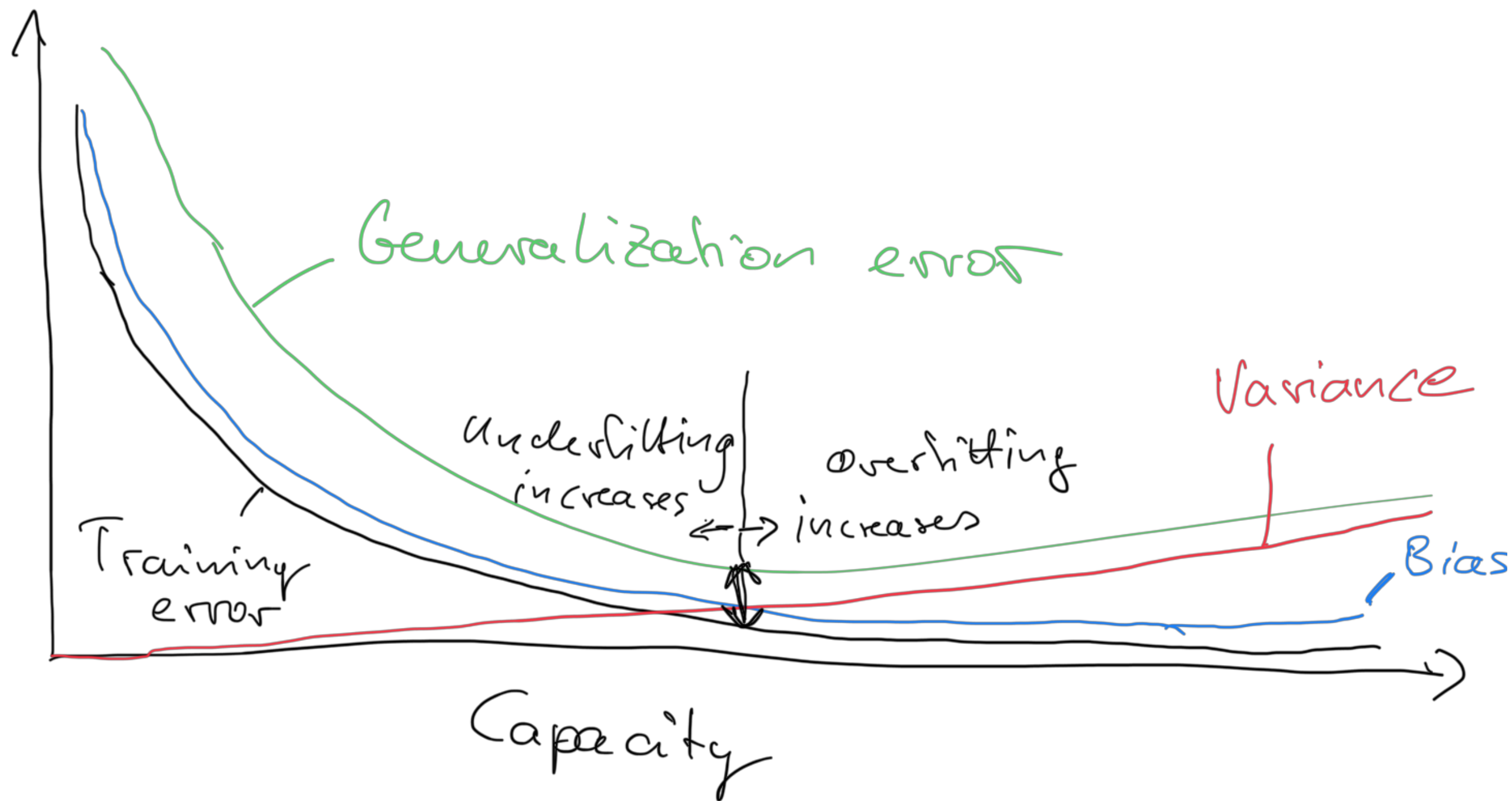$$= \textbf{[Bias of the fit]}^2 + \textbf{Variance of the fit}$$

# Bias-Variance Decomposition of Squared Error

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$\text{???}$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - y)^2 \boxed{+ 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})}$$

$$
\begin{aligned}
E[2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] &= 2E[(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] \\
&= 2(y - E[\hat{y}])E[(E[\hat{y}] - \hat{y})] \\
&= 2(y - E[\hat{y}])(E[E[\hat{y}]] - E[\hat{y}]) \\
&= 2(y - E[\hat{y}])(E[\hat{y}] - E[\hat{y}]) \\
&= 0
\end{aligned}
$$

Generalization error

Variance

Underfitting
increases

Overfitting
increases

Training
error

Bias

Capacity

# to be continued ...