

# Overfitting and Underfitting

These slides are based on the slides of Sebastian rashcka for introduction to deep learning course and the slides of Sharifi Zarchi for the introduction to machine learning course

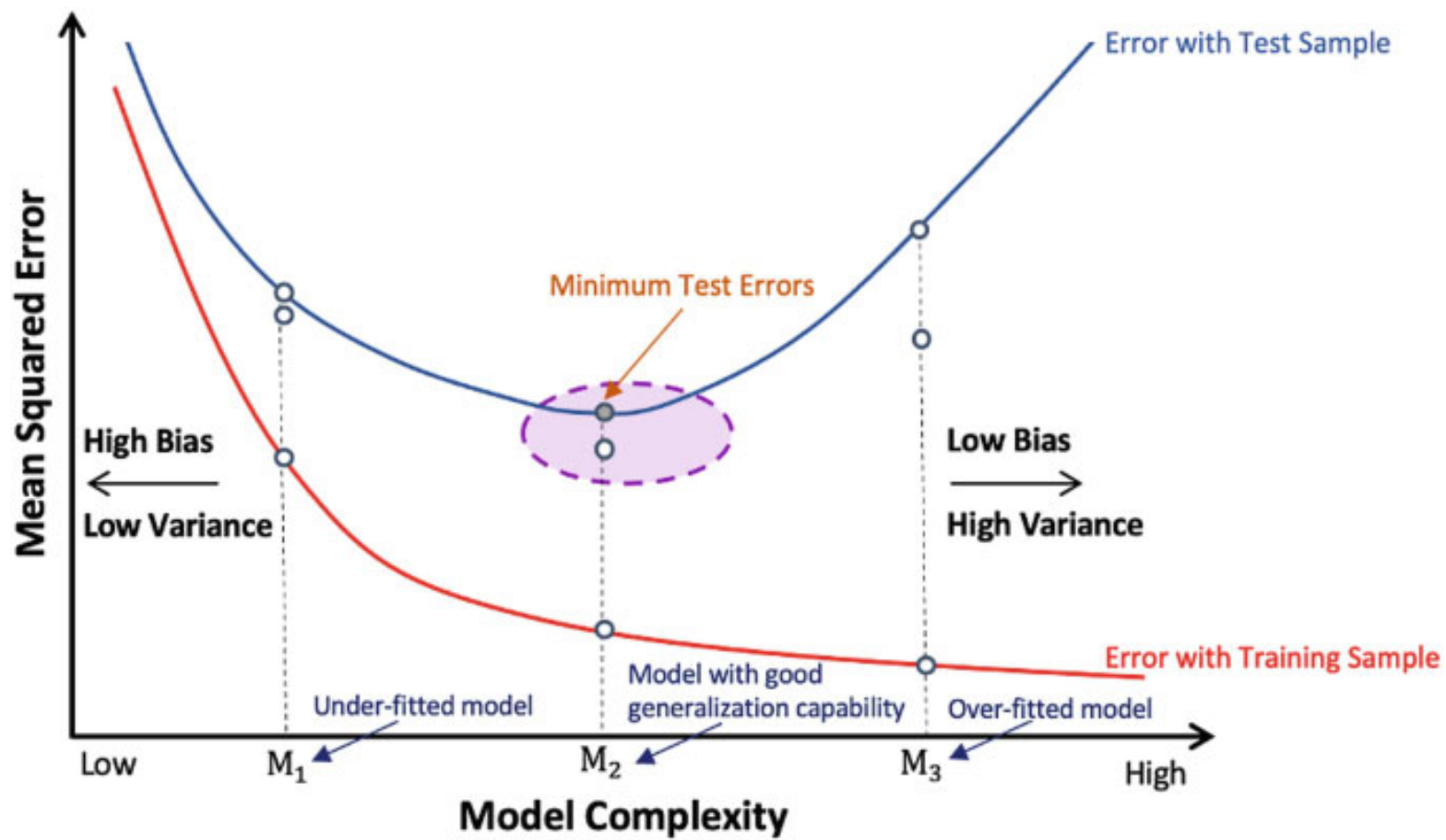
# Generalization Performance

- Want a model to "generalize" well to \_\_\_\_\_ data
- (Want "high generalization accuracy" or "low generalization error")
- Test sets is a method to estimating the generalization performance on unseen data.

# Assumptions

- i.i.d. assumption: training and test examples are independent and identically distributed (drawn from the same joint probability distribution,  $P(X, y)$  )
- For some random model that has not been fitted to the training set, we expect **the training error is \_\_\_\_\_ the test error**

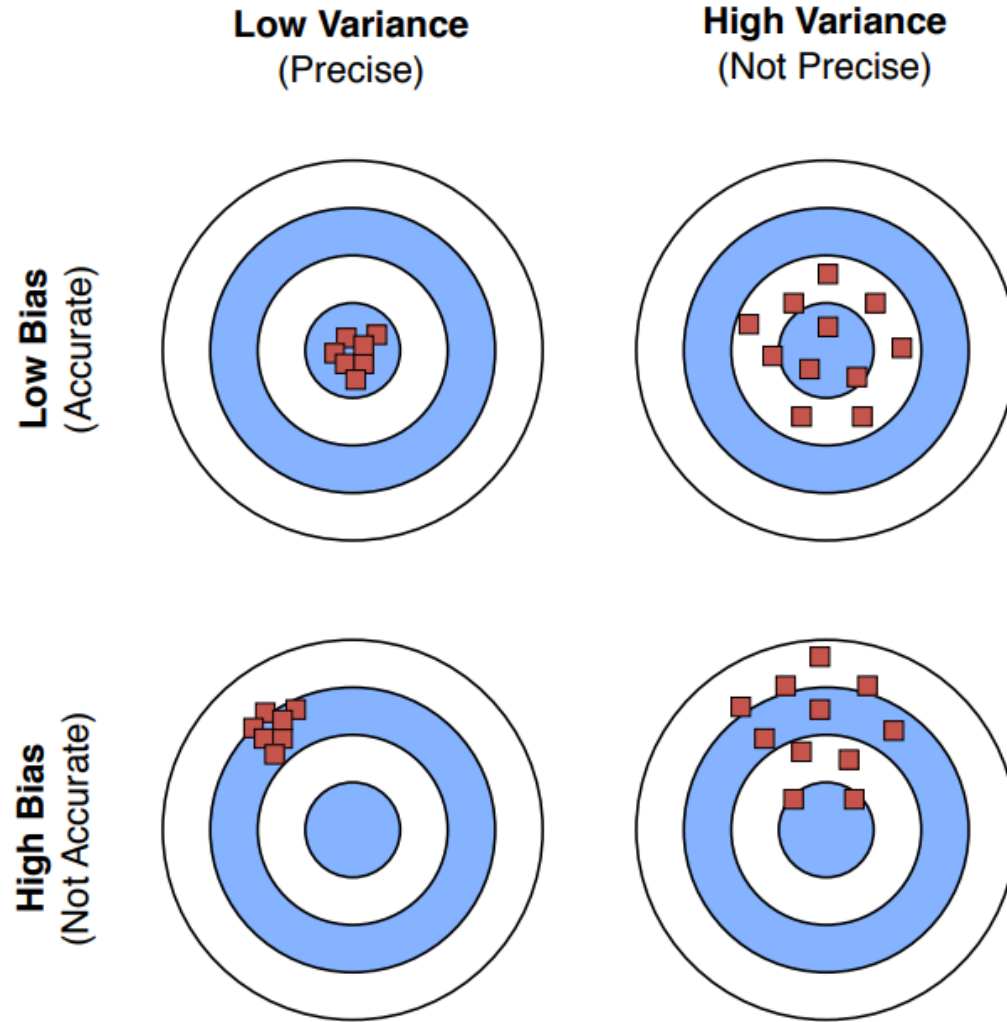
- Underfitting: both the training and test error are \_\_\_\_\_
- Overfitting: gap between training and test error (where test error is larger)
- Large hypothesis space being searched by a learning algorithm -> high tendency to \_\_\_\_\_ fit



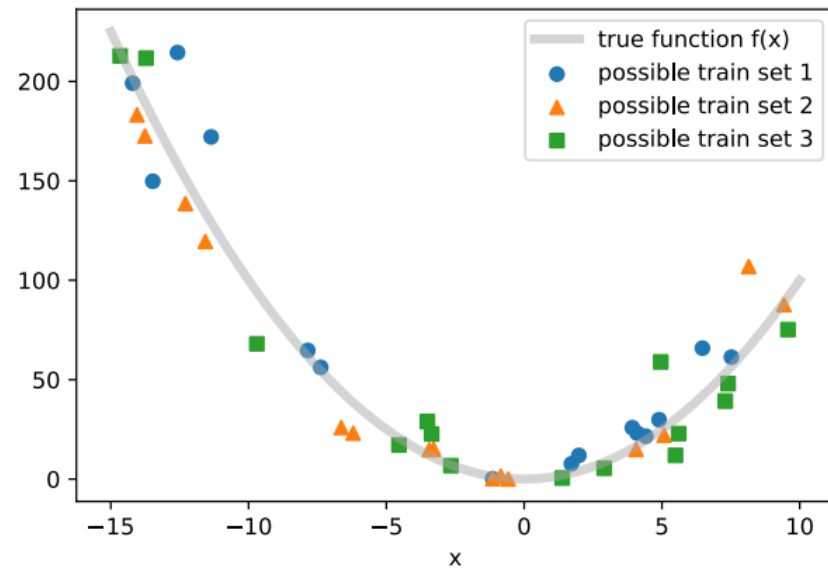
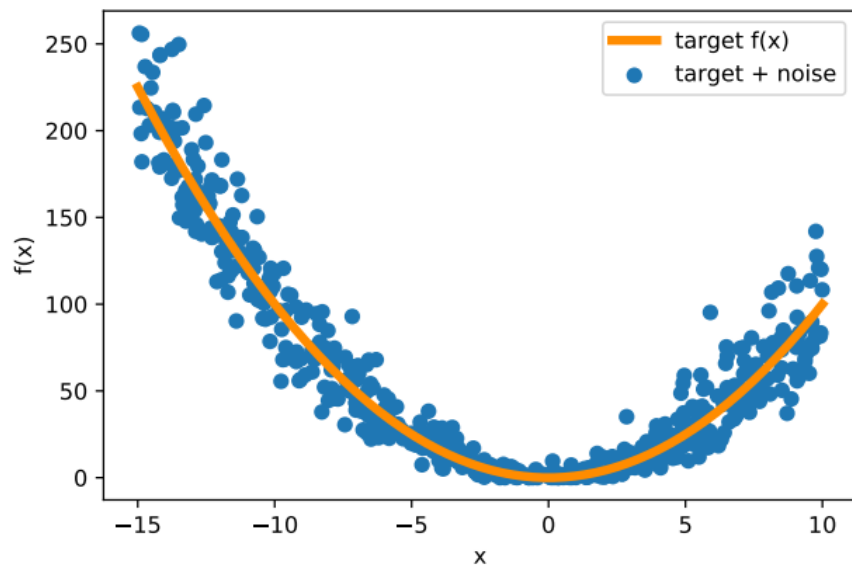
# Bios – Variance decomposition

- Decomposition of the loss into bias and variance help us understand learning algorithms, concepts are related to underfitting and overfitting
- Helps explain why ensemble methods might perform better than single models

# Bias-Variance Intuition

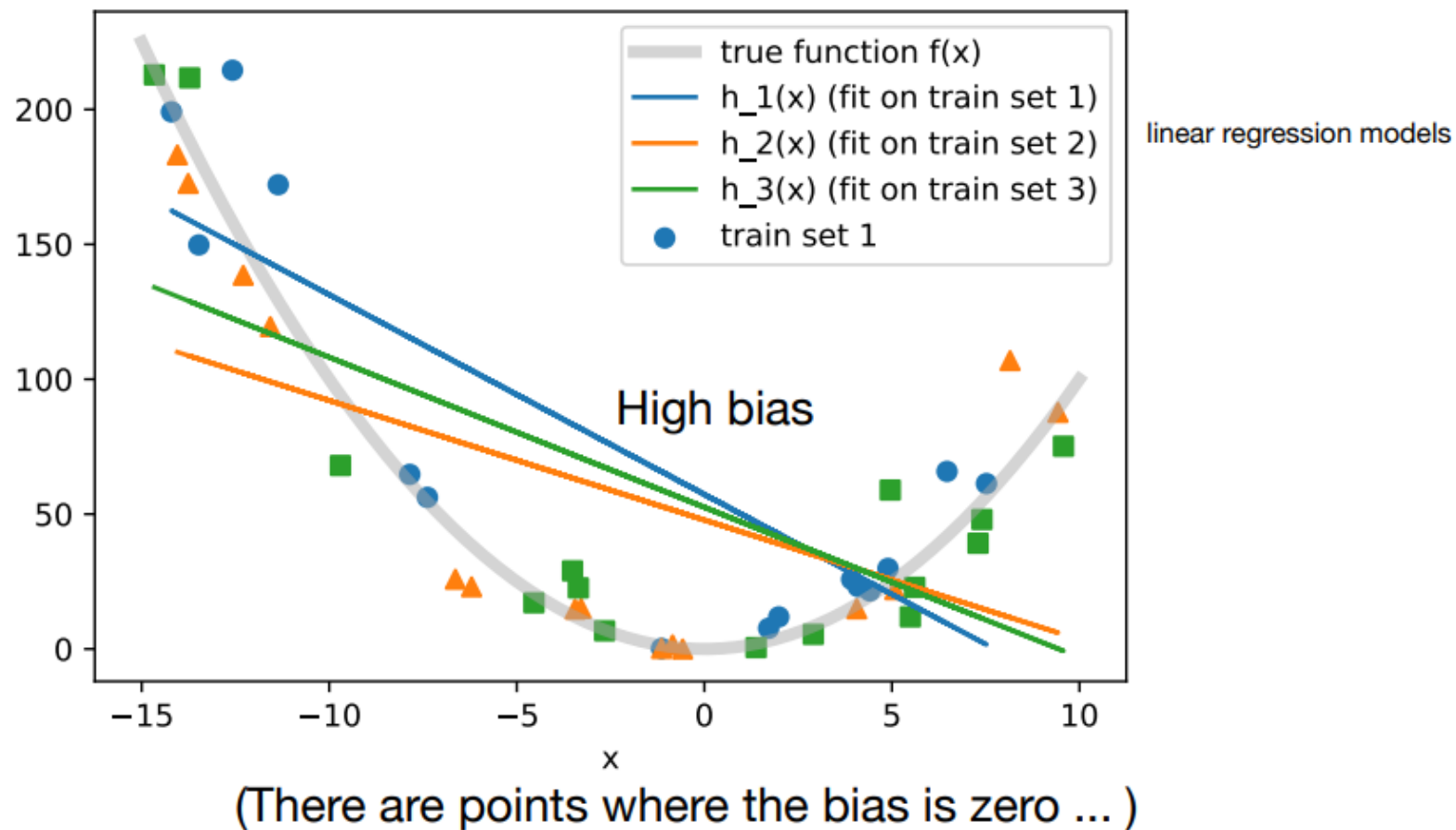


# Bias-Variance Intuition



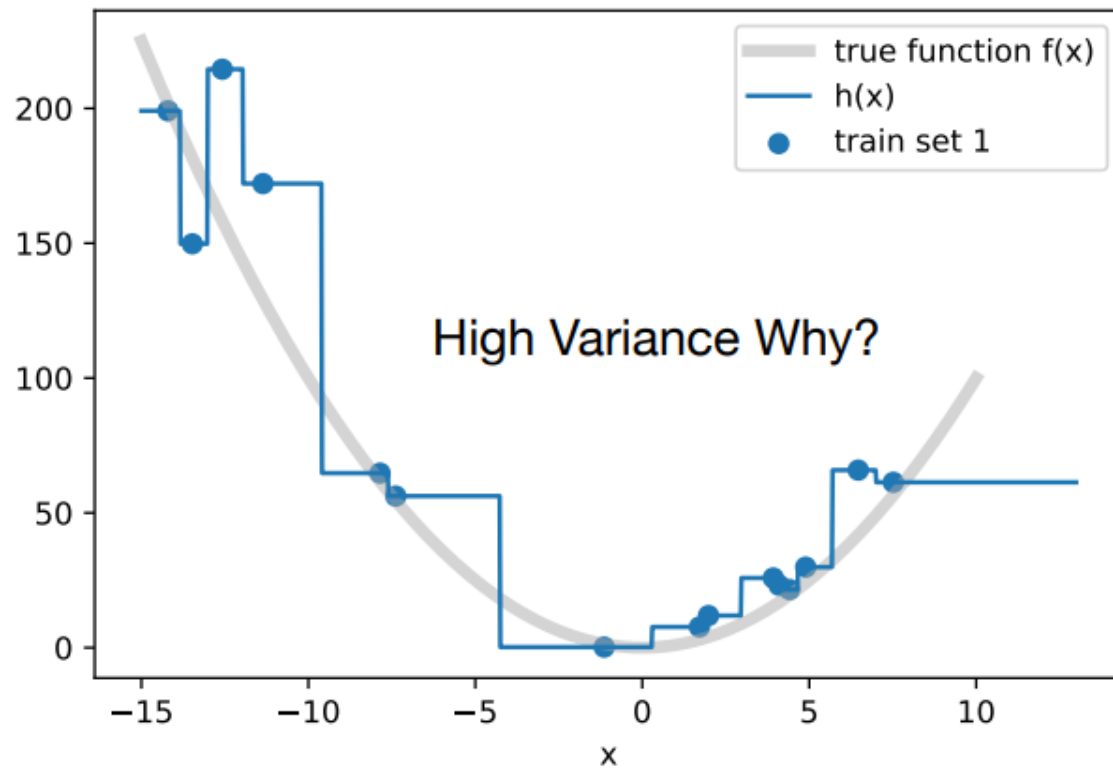


# Bias-Variance Intuition



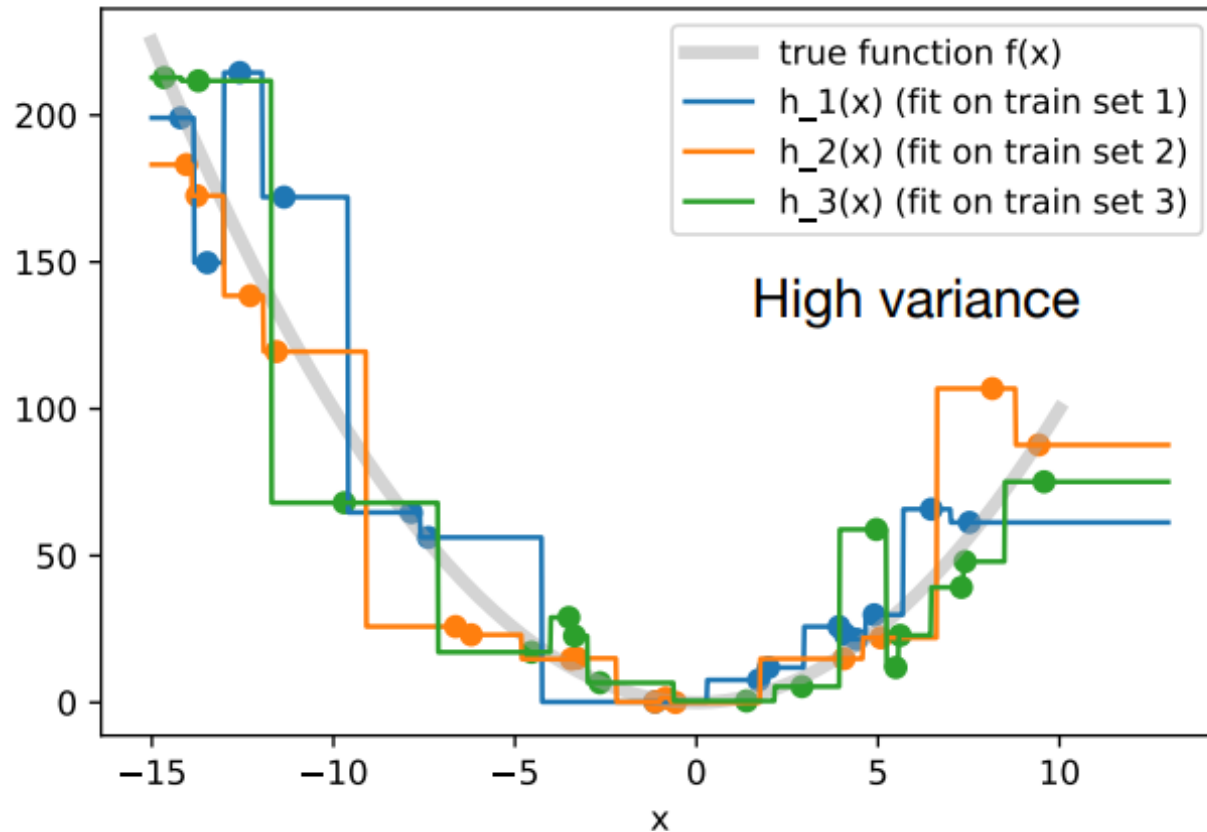
# Bias-Variance Intuition

(here, I fit an unpruned decision tree)

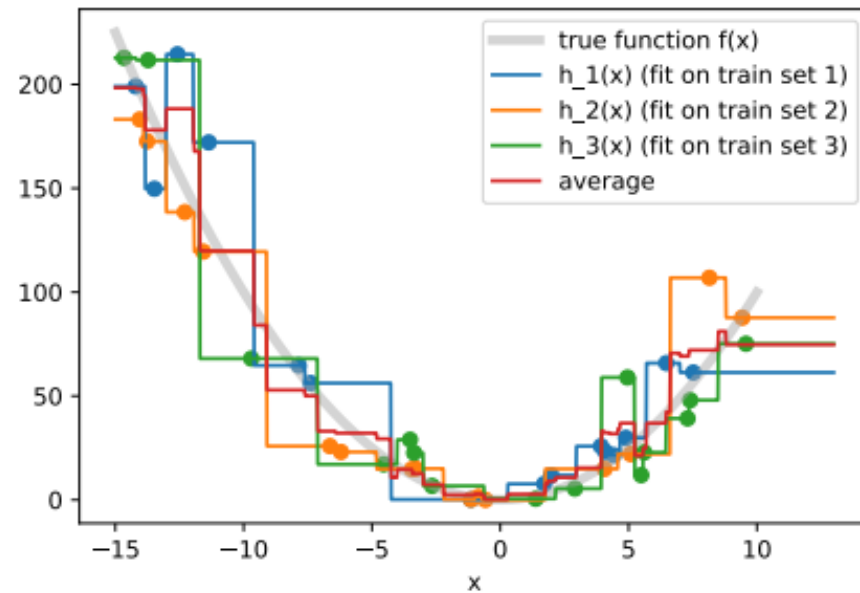
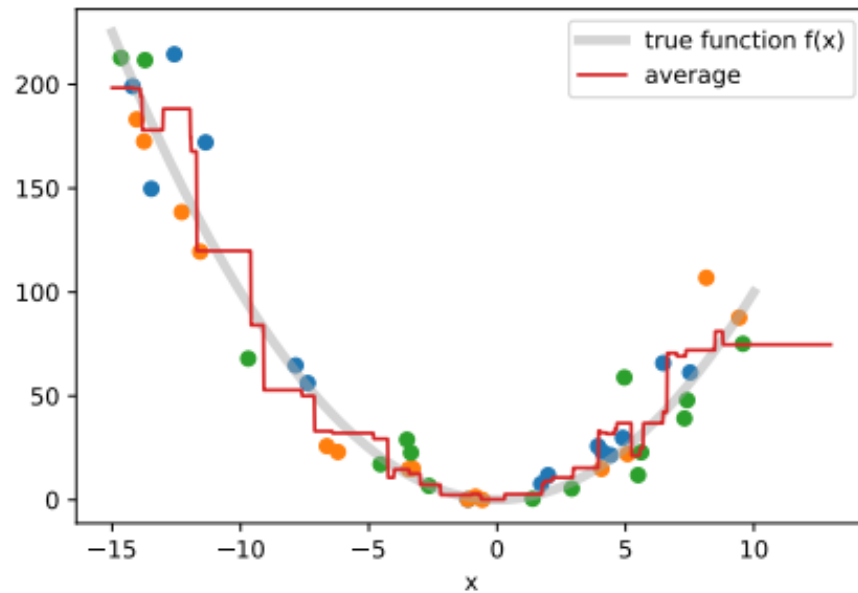


# Bias-Variance Intuition

suppose we have multiple training sets



# Bias-Variance Intuition



# Terminology

Point estimator  $\hat{\theta}$  of some parameter  $\theta$

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

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

# Terminology

## General Definition

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

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

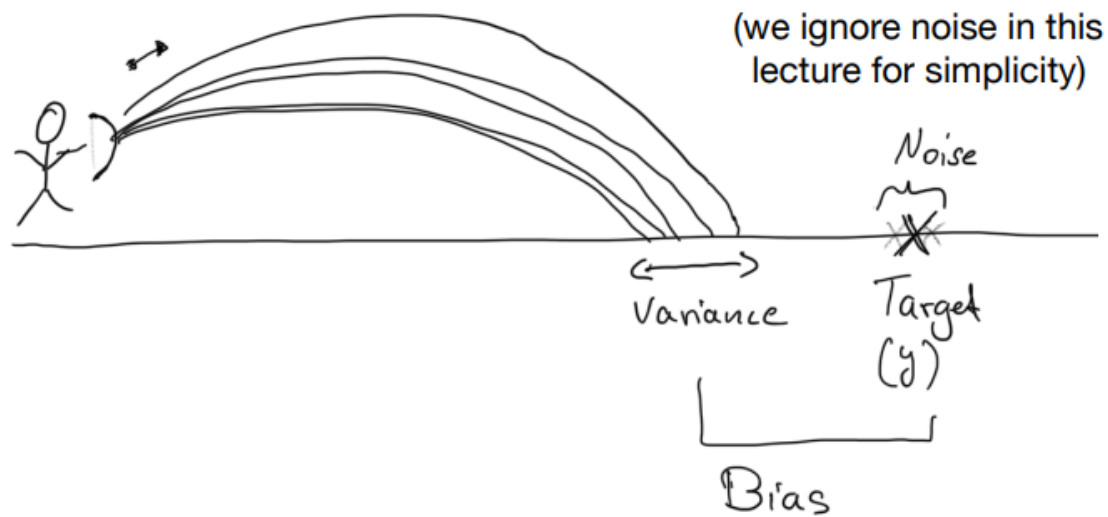
$$\text{Var}[\hat{\theta}] = E \left[ (E[\hat{\theta}] - \hat{\theta})^2 \right]$$

# Terminology

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

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

## Intuition



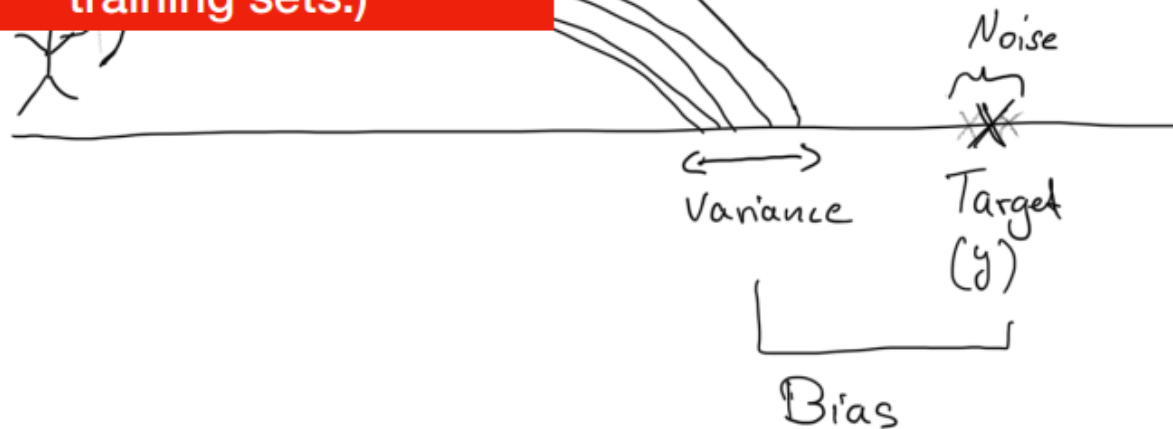
# Terminology

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

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

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

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





# Bias-Variance of the Squared Error

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

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

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


## "ML Notation" for Squared Error Loss

$y = f(x)$  target

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

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

for simplicity, we ignore  
the noise term



# Bias-Variance of the Squared Error

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

**"ML Notation" for  
Squared Error Loss**

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

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

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

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

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

## Bias-Variance of the Squared Error

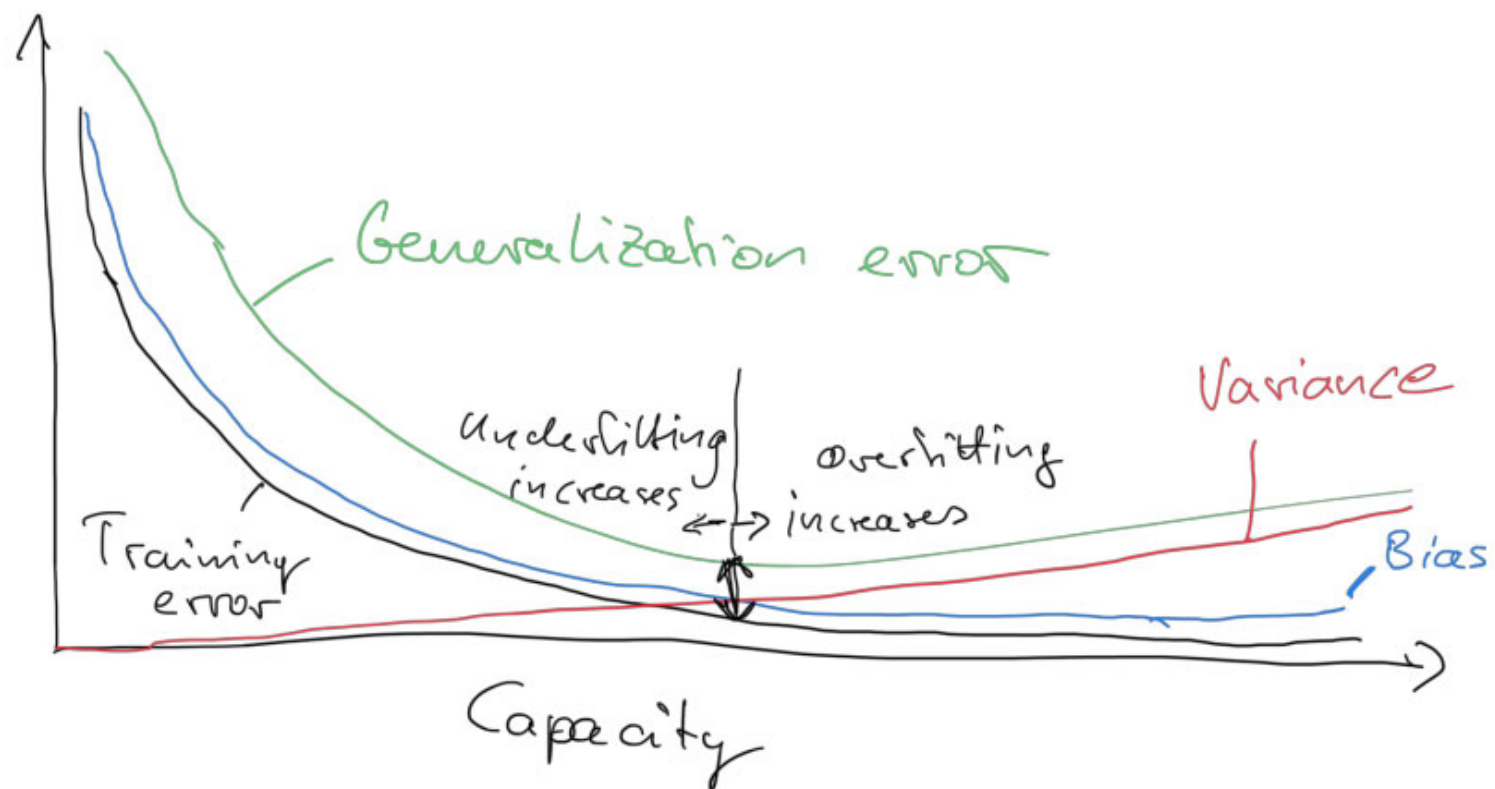
$$S = (y - \hat{y})^2$$

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

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

$$\begin{aligned}E[(y - \hat{y})^2] &= (y - E[\hat{y}])^2 + E[(E[\hat{y}] - \hat{y})^2] \\ &= \text{Bias}^2 + \text{Var}\end{aligned}$$

**Now, how is this related to overfitting and underfitting?**



# Overfitting

- **Overfitting** means the model works well on training data, but it doesn't generalize well.
- Overfitting occurs when there is too much complexity in the model in comparison to the amount and noise in the training data.

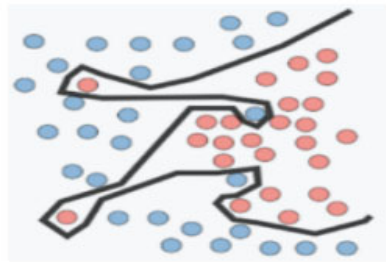


Figure: High Variance, [Source](#)

- How to fix this problem?
  - ▷ Simplify the model by selecting one with fewer parameters, reducing the number of attributes in the training data, or constraining the model.
  - ▷ Gather more training data.
  - ▷ Reduce the noise in the training data.

# Underfitting

- **Underfitting** is the opposite of overfitting: it occurs when your model is too simple to learn the underlying structure of the data.

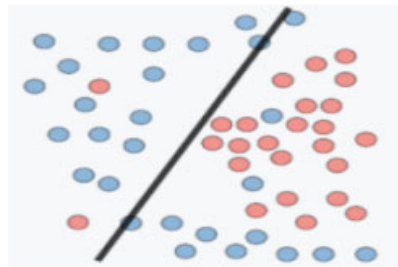


Figure: High Variance, [Source](#)

- The main options to fix this problem:
  - ▷ Selecting a more powerful model, with more parameters
  - ▷ Feeding better features to the learning algorithm (feature engineering)
  - ▷ Reducing the constraints on the model (e.g., reducing the regularization hyperparameter)