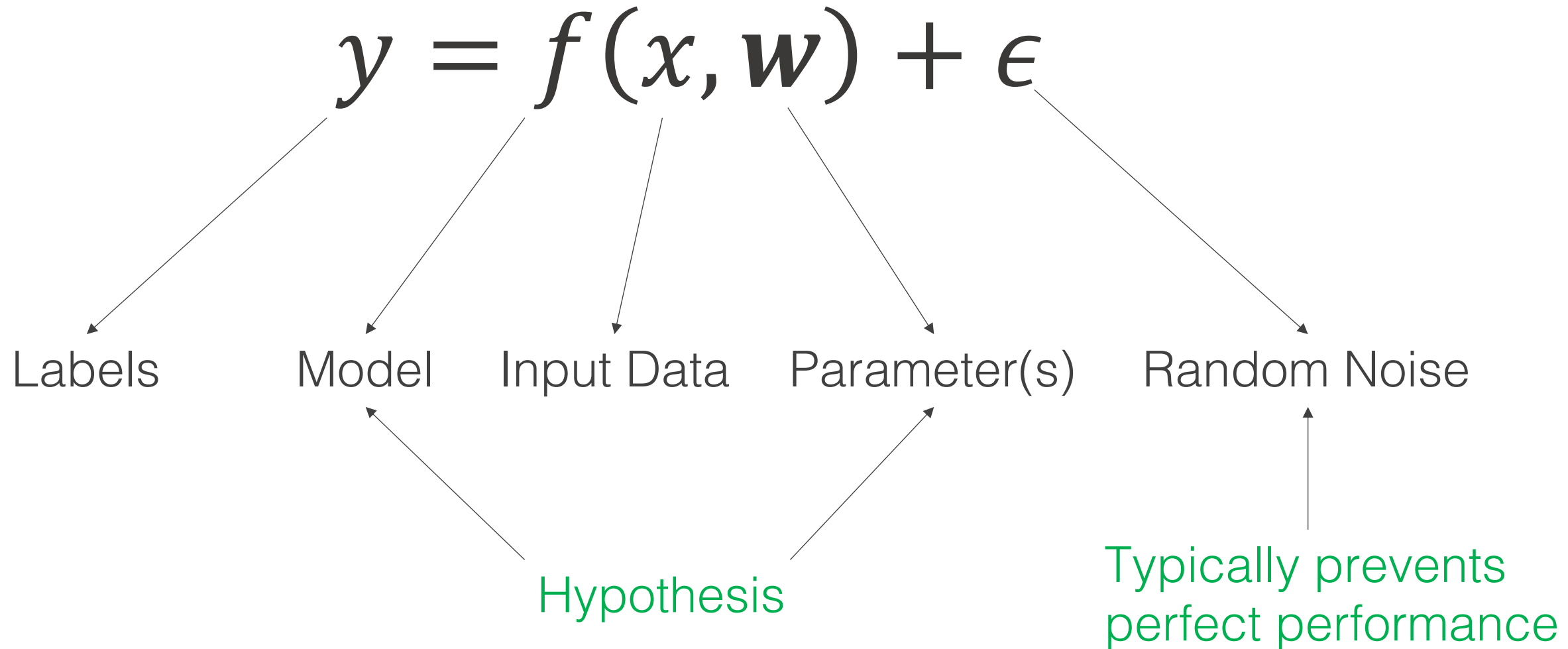# How flexible should my algorithms be?

Lecture 03

# Supervised Learning

Algorithm development and application pipeline

# Supervised machine learning model
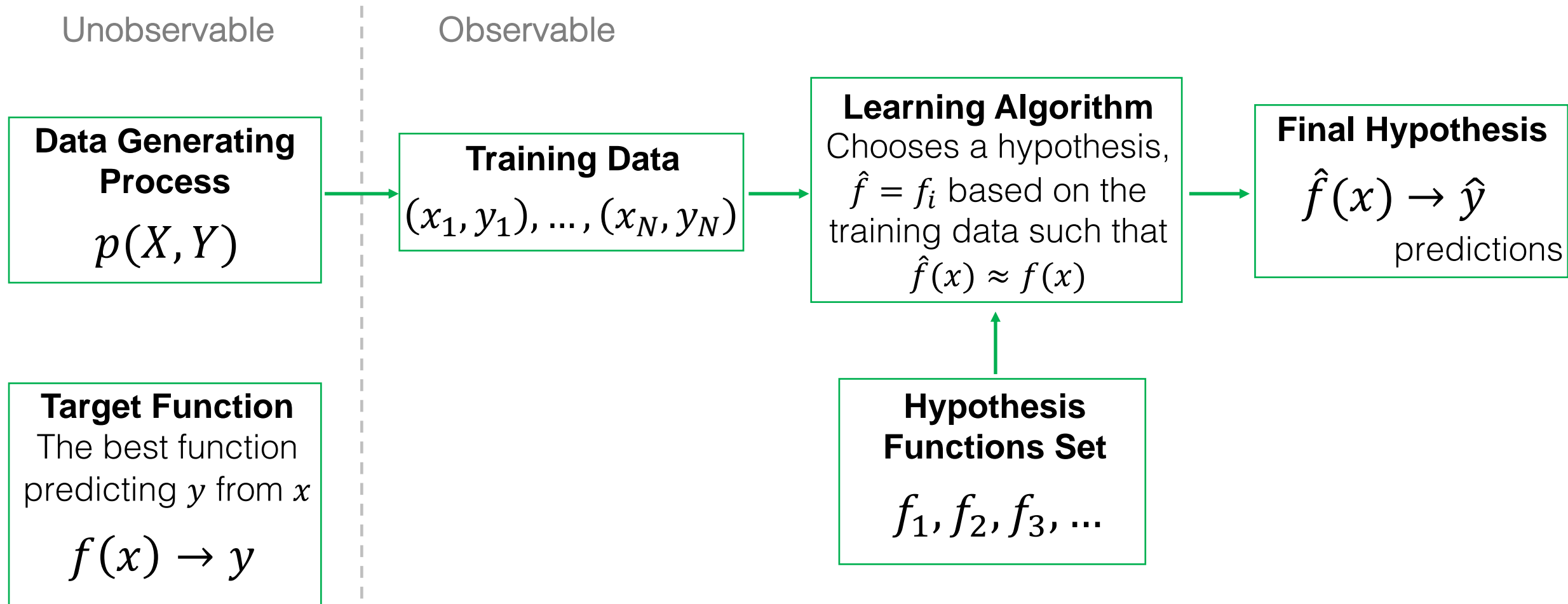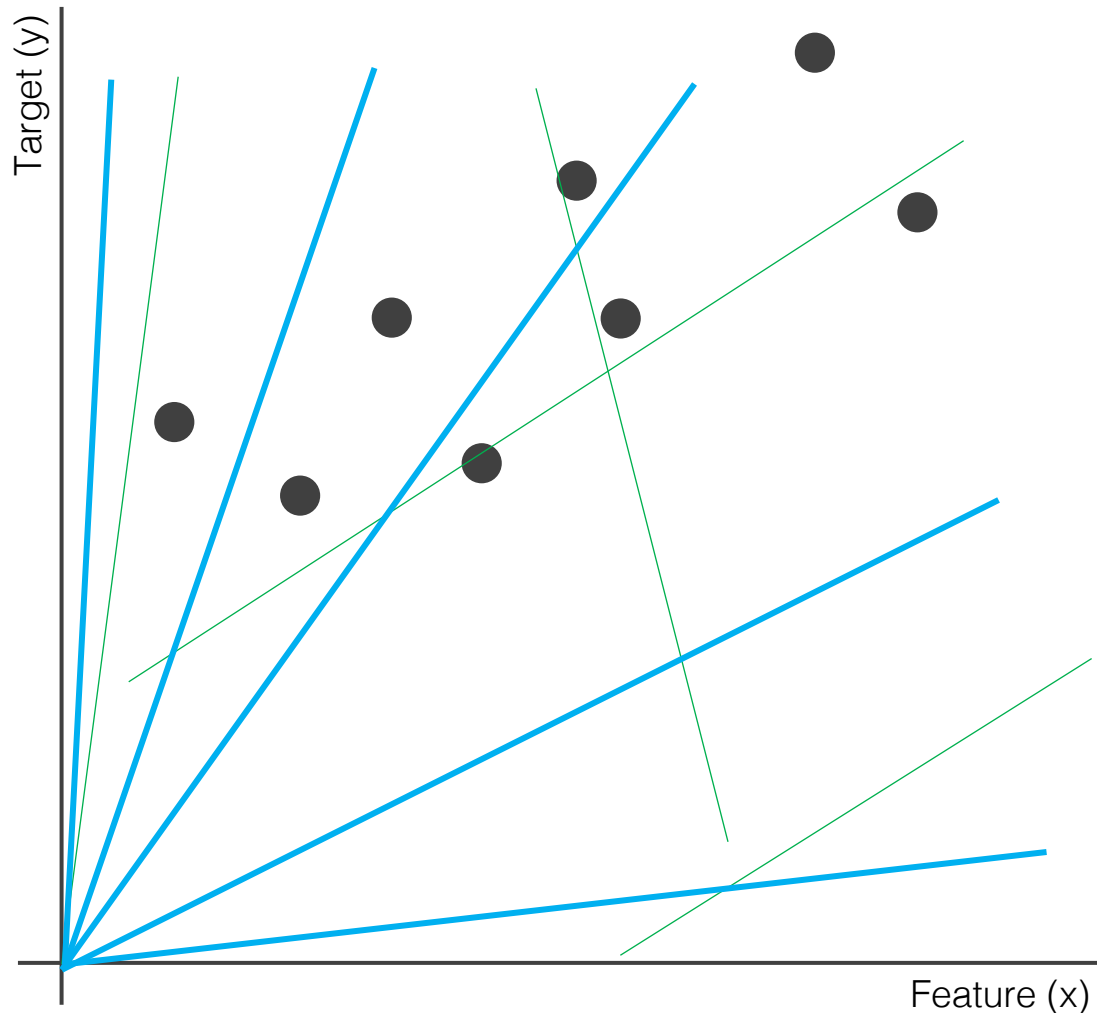
We search for the model that best fits our data

$$y = f(x, \boldsymbol{w}) + \epsilon$$

Labels    Model    Input Data    Parameter(s)    Random Noise

Hypothesis

Typically prevents perfect performance

# Components of **supervised** learning

**Input** $\qquad\qquad\qquad x$

**Output** $\qquad\qquad\quad y$

**Training Data** $\qquad (x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

**Target function** $\qquad f(x) \rightarrow y \qquad$ This is unknown, but the best you could ever do

**Hypothesis set** $\qquad f_i(x) \rightarrow \hat{y} \qquad$ Functions to consider in trying to approximate $f(x)$

**Learning algorithm** $\quad$ Optimization technique that searches the hypothesis set for the function $f_i$ that best approximates $f$ (typically by choosing parameters in a model)

# **Supervised** Learning

**Data Generating Process**

$$p(X, Y)$$

**Training Data**

$$(x_1, y_1), \ldots, (x_N, y_N)$$

**Learning Algorithm**
Chooses a hypothesis, $\hat{f} = f_i$ based on the training data such that $\hat{f}(x) \approx f(x)$

**Final Hypothesis**

$$\hat{f}(x) \rightarrow \hat{y}$$
predictions

**Target Function**
The best function predicting $y$ from $x$

$$f(x) \rightarrow y$$

**Hypothesis Functions Set**

$$f_1, f_2, f_3, \ldots$$

- Need to select the hypothesis functions (models to train)
- Need to select the learning algorithm (for fitting the models to the data)

# Example: linear regression



Using any line as a hypothesis function, how many possible hypothesis functions are in the set?

**Infinitely many**

Using the line $y = wx$ as the family of hypothesis functions, how many possible hypothesis functions are in the set?
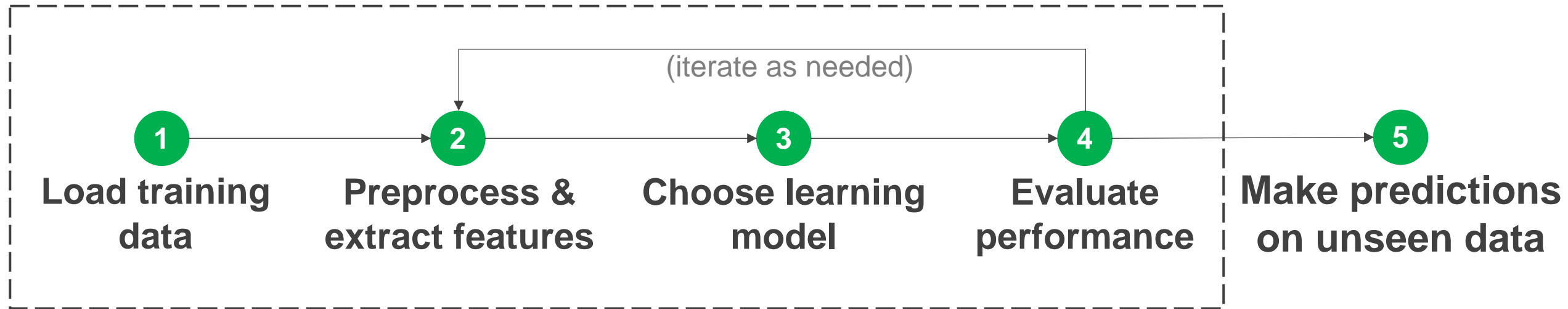
**Infinitely many**

Which set contains the better hypothesis?

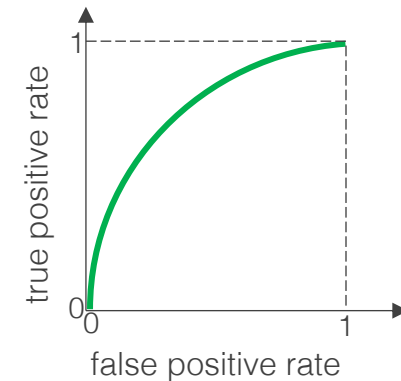Which set has more options to consider?

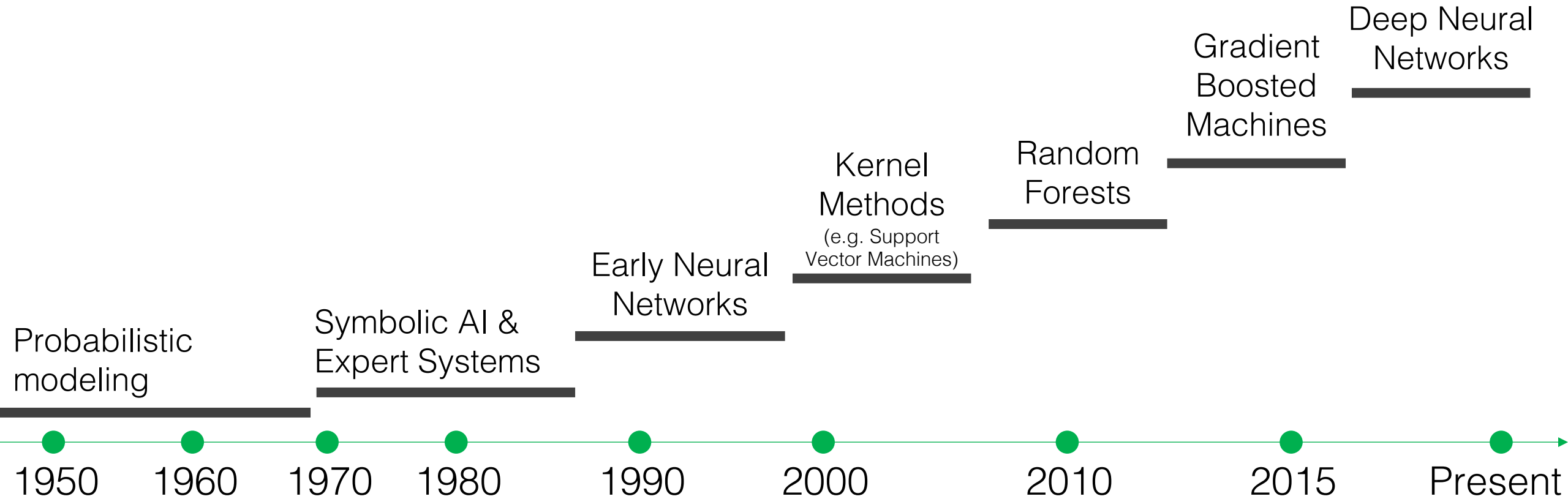What is our learning algorithm?

# Algorithm Development

## Application

(iterate as needed)

**1** Load training data → **2** Preprocess & extract features → **3** Choose learning model → **4** Evaluate performance → **5** Make predictions on unseen data

$y$

| | |
|---|---|
| 1 | |
| 1 | |
| 0 | |
| 0 | |
| 1 | |
| 0 | |

**X**

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$
$x_6$

**X'**

| | | | | |
|---|---|---|---|---|
| $x_1$ | 0.38 | 0.39 | 0.85 | 0.78 |
| $x_2$ | 0.81 | 0.91 | 0.97 | 0.53 |
| $x_3$ | 0.65 | 0.59 | 0.91 | 0.11 |
| $x_4$ | 0.94 | 0.05 | 0.40 | 0.26 |
| $x_5$ | 0.27 | 0.19 | 0.03 | 0.64 |
| $x_6$ | 0.02 | 0.98 | 0.36 | 0.11 |

linear discriminant
perceptron
logistic regression
decision trees
random forests
support vector machine
k nearest neighbors
neural networks

true positive rate / false positive rate

$X_{new}$   $\widehat{y}_{new}$

→ **1**

→ **0**

# Historic Progression of Algorithms



Deep Neural Networks

Gradient Boosted Machines

Random Forests

Kernel Methods
(e.g. Support Vector Machines)

Early Neural Networks

Symbolic AI & Expert Systems

Probabilistic modeling

| 1950 | 1960 | 1970 | 1980 | 1990 | 2000 | 2010 | 2015 | Present |

François Chollet, *Deep Learning with Python*, 2017

# How flexible should my algorithms be?

Lecture 03

# K-Nearest Neighbors
## Classification and Regression

# K Nearest Neighbor Classifier

**Step 1**:
Training

Every new data point is a model parameter

# K Nearest Neighbor Classifier

**Step 2**: Place new (unseen) examples in the feature space

# K Nearest Neighbor Classifier

**Step 3**: Classify the data by assigning the class of the k nearest neighbors

# K Nearest Neighbor Classifier



**Score vs Decision** :

For 5-NN, the confidence score that a sample belongs to a class could be: $\{0, 1/5, 2/5, 3/5, 4/5, 1\}$

**Decision Rule**:

If the confidence score for a class > threshold, predict that class

# K Nearest Neighbor Regression

# K Nearest Neighbor Regression

$$\hat{y} = \frac{1}{k} \sum_{y_i \in \{\text{k nearest}\}} y_i$$



$\hat{y} \cong 3.67$

# KNN Pros and Cons

## Pros

- Simple to implement and interpret
- Minimal training time
- Naturally handles multiclass data

## Cons

- Computational expensive to find nearest  neighbors
- Requires all of the training data to be stored in the model
- Suffers if classes are imbalanced
- Performance may suffer in high dimensions

# How flexible should my model be?

the bias-variance tradeoff and learning to generalize

**bias**
consistently incorrect prediction

error from poor model assumptions
(high bias results in underfit)

**variance**
inconsistent prediction

error from sensitivity to small changes in the training data
(high variance results in overfit)

**noise**
lower bound on generalization error

irreducible error inherent to the problem
(e.g. you cannot predict the outcome of a flip of a fair coin any more than 50% of the time)

# Bias-Variance Tradeoff

generalization error  =  **bias**$^2$  +  **variance**  +  noise

# Classification feature space



How flexible should my algorithms be? Duke University | Lecture 03 21

# What's the best we can do for binary classification?

If we know the probability distribution of the data

The Bayes decision rule

# Bayes' Rule

$$P(C|X) = \frac{P(X|C)P(C)}{P(X)}$$

Likelihood   Prior

Posterior

Evidence

$X$   Features

$C$   Class label

i.e. C $\in \{c_0, c_1\}$ for the binary case

## Bayes' Decision Rule:
choose the most probable class given the data

If $P(C_i = c_1|X_i) > P(C_i = c_0|X_i)$ then $\hat{y} = c_1$

otherwise $\hat{y} = c_0$

- If the distributions are correct, this decision rule is **optimal**

- Rarely do we have enough information to use this in practice

# Classification feature space

Bayes Decision Boundary

# Decision Boundary Examples



Bayes Decision Boundary

Linear Classifier

Decision Tree

k=1 Nearest Neighbor

k=15 Nearest Neighbor

Support Vector Machine
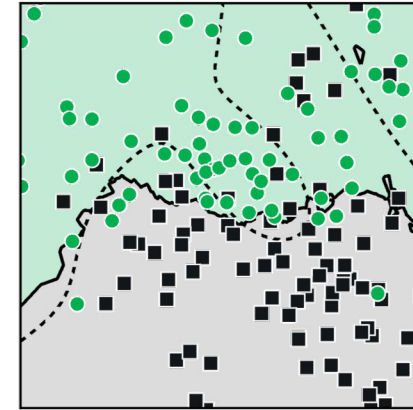
# Bias Variance Tradeoff

k = 1    k = 15    k = 100
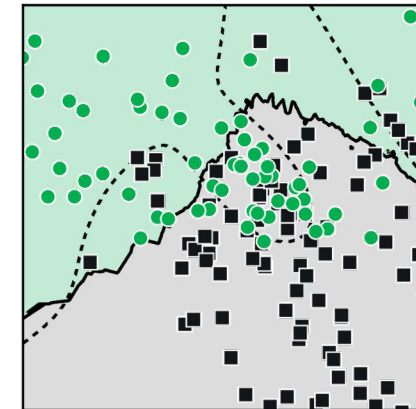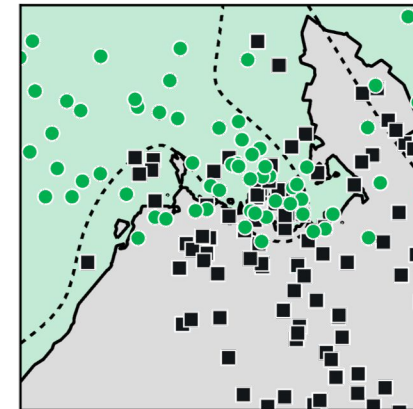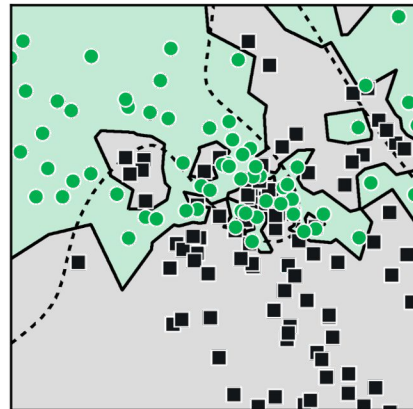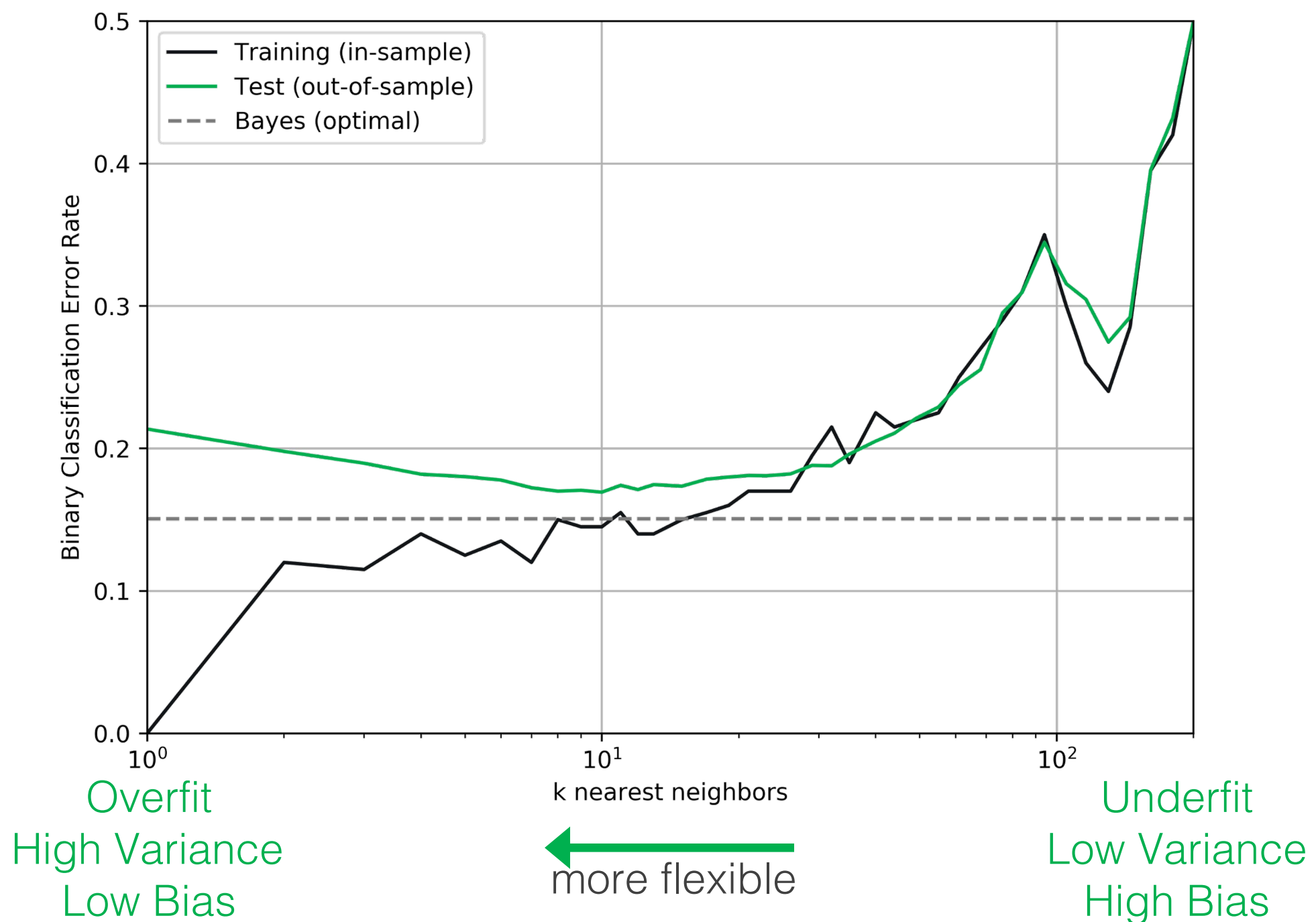
Sample 1

higher bias

underfit

Sample 2

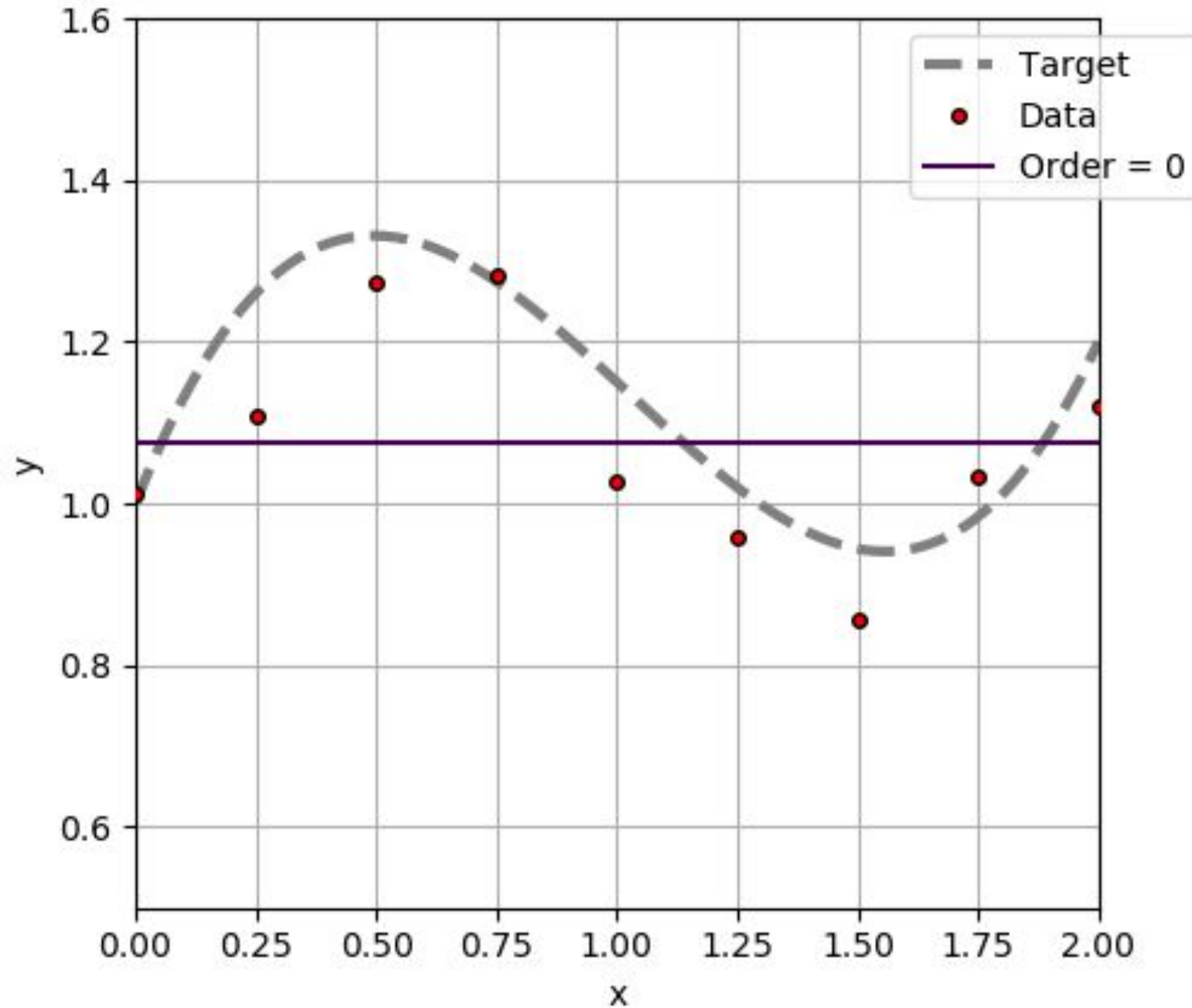higher variance

overfit

Sample 3

# Bias Variance Tradeoff

# This tradeoff is equally challenging for regression

# Linear Regression

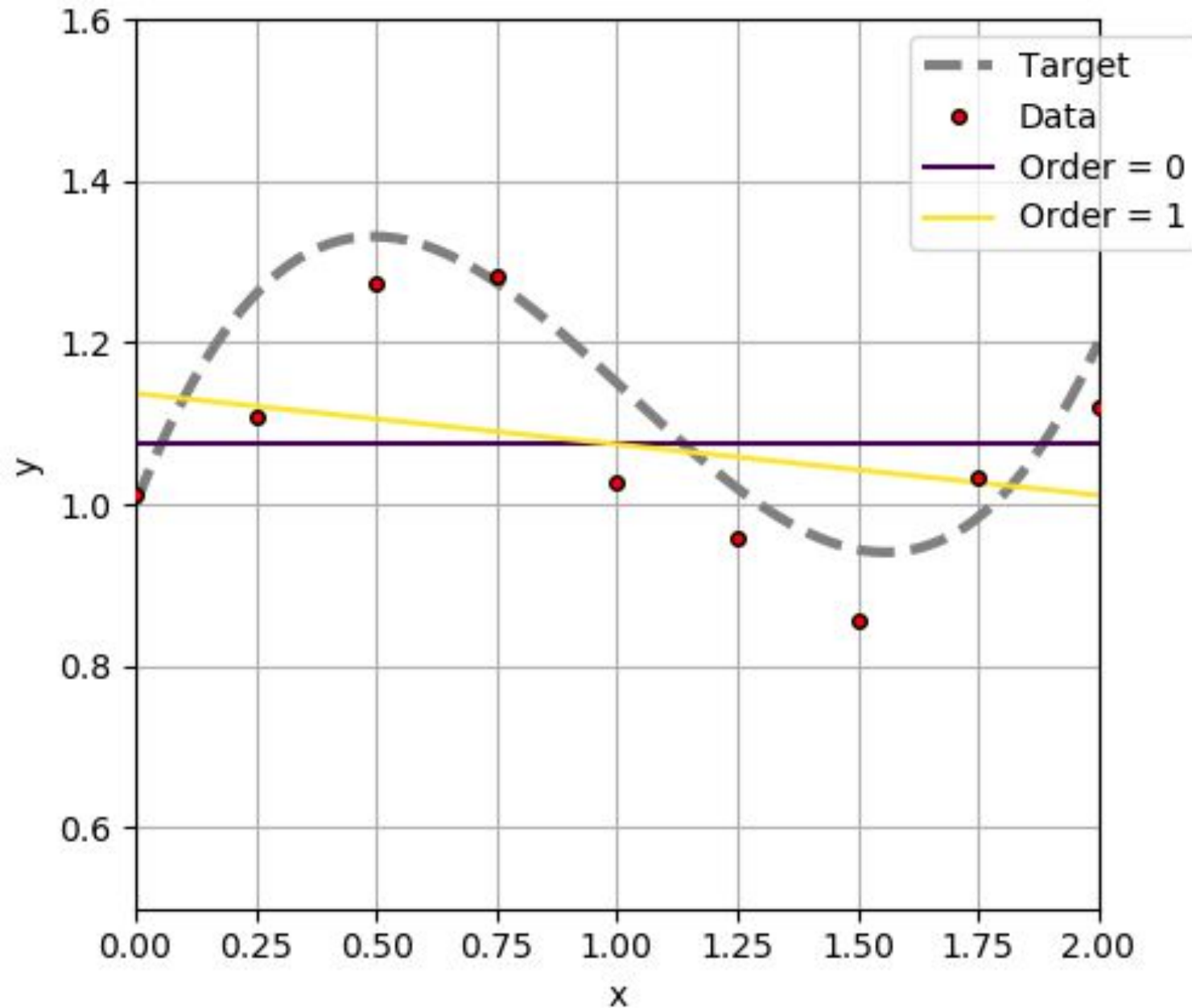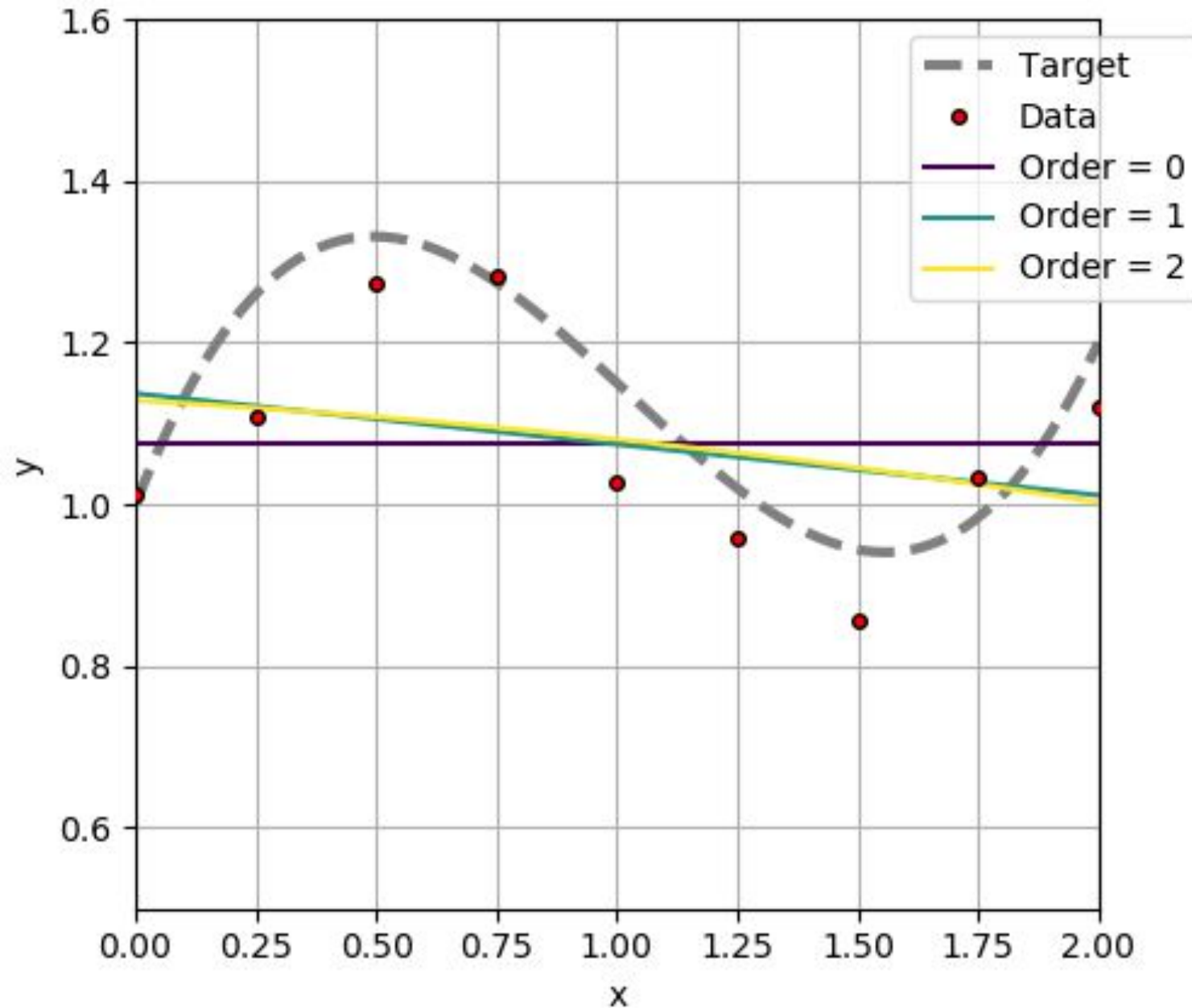$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

m is the model order

# Linear Regression

$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

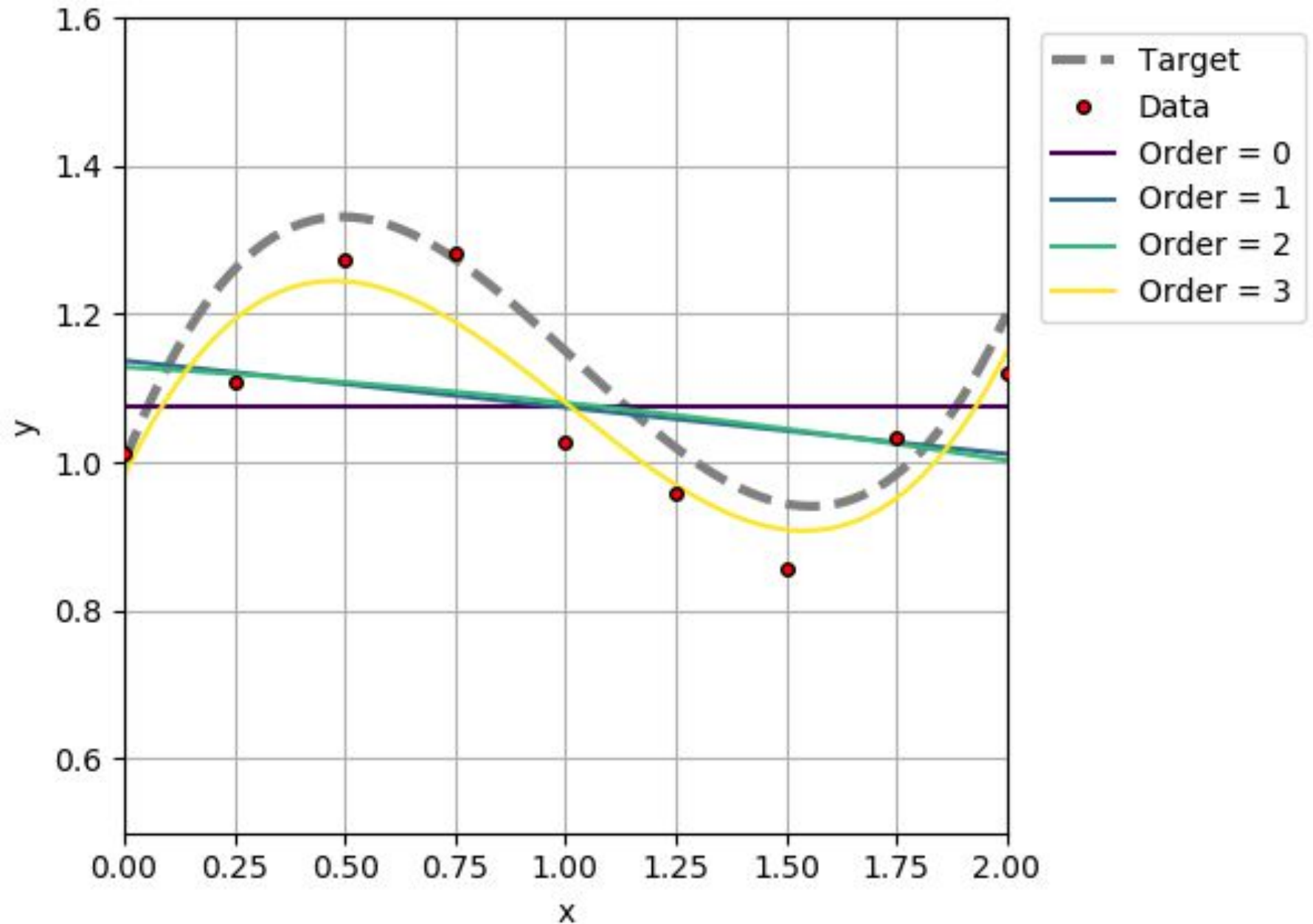m is the model order

# Linear Regression

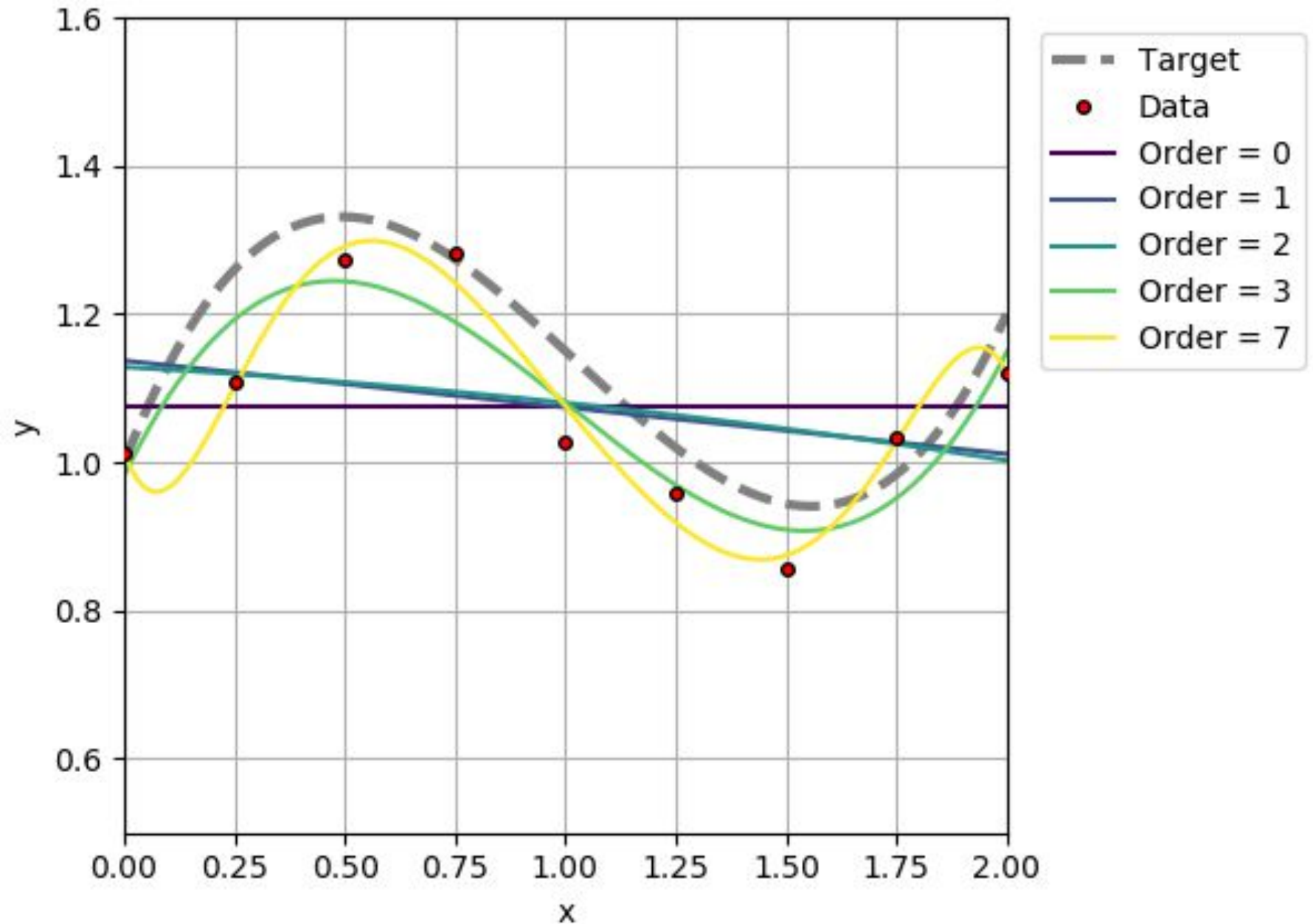$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

m is the model order

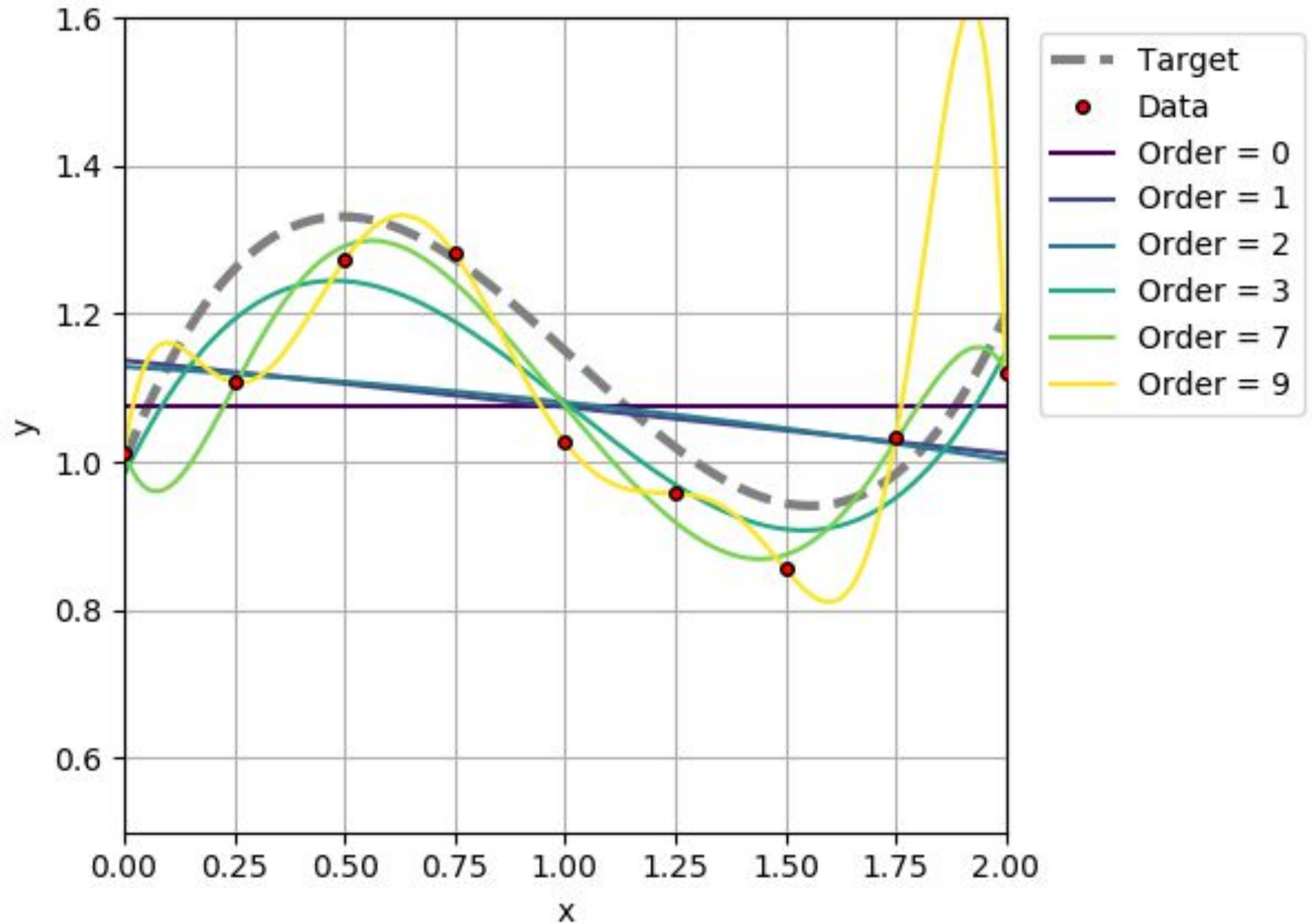# Linear Regression

$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

m is the model order

# Linear Regression

$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

m is the model order



Legend:
- Target
- Data
- Order = 0
- Order = 1
- Order = 2
- Order = 3
- Order = 7

# Linear Regression

$$\hat{y}_i = \sum_{j=0}^{m} a_j x_i^j$$

m is the model order

# Problem

Too much flexibility leads to **overfit**

Too little flexibility leads to **underfit**

Over/underfit **hurts generalization** performance

# Solutions for overfitting

1. Add **more data** for training

2. Constrain model flexibility through **regularization**

3. Use model **ensembles**