

How flexible should my algorithms be?

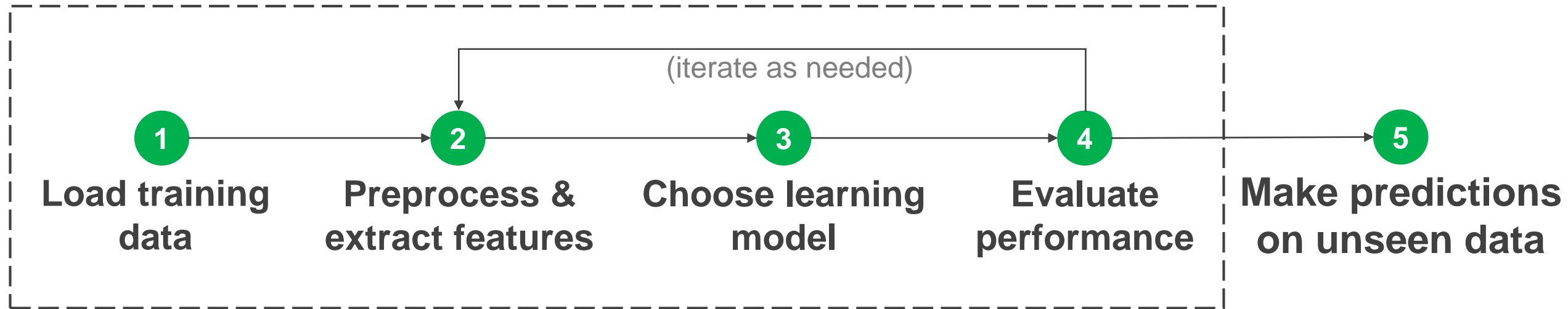
Lecture 03

Review of Supervised Learning

Algorithm development and application pipeline

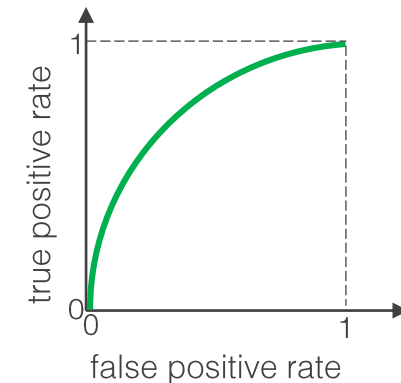
Algorithm Development

Application



y	X	X'
1	x_1	x_1 0.38 0.39 0.85 0.78
1	x_2	x_2 0.81 0.91 0.97 0.53
0	x_3	x_3 0.65 0.59 0.91 0.11
0	x_4	x_4 0.94 0.05 0.40 0.26
1	x_5	x_5 0.27 0.19 0.03 0.64
0	x_6	x_6 0.02 0.98 0.36 0.11

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



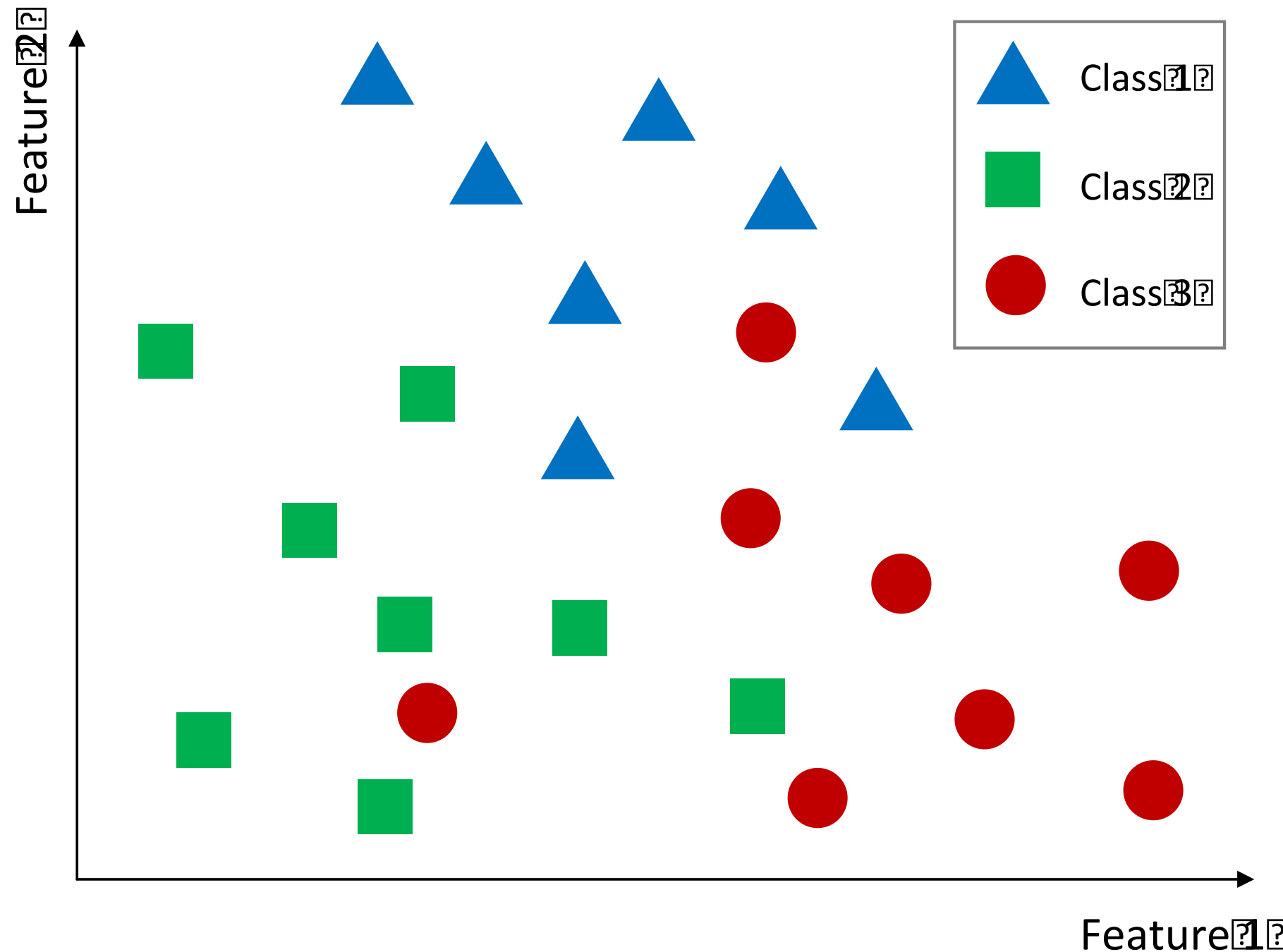
K-Nearest Neighbors

Classification and Regression

K Nearest Neighbor Classifier

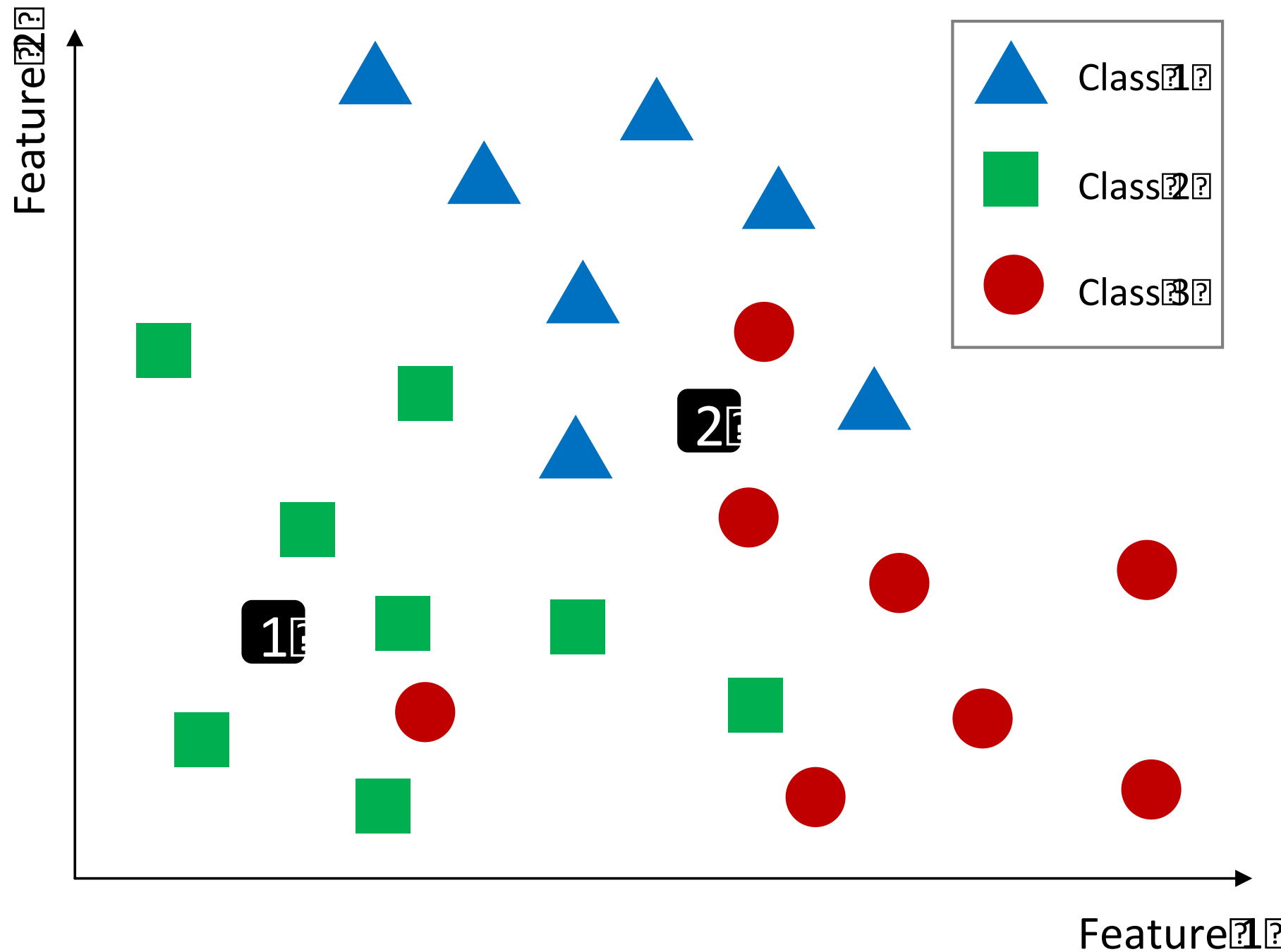
Step 1: Training

Every new data point is
a model parameter



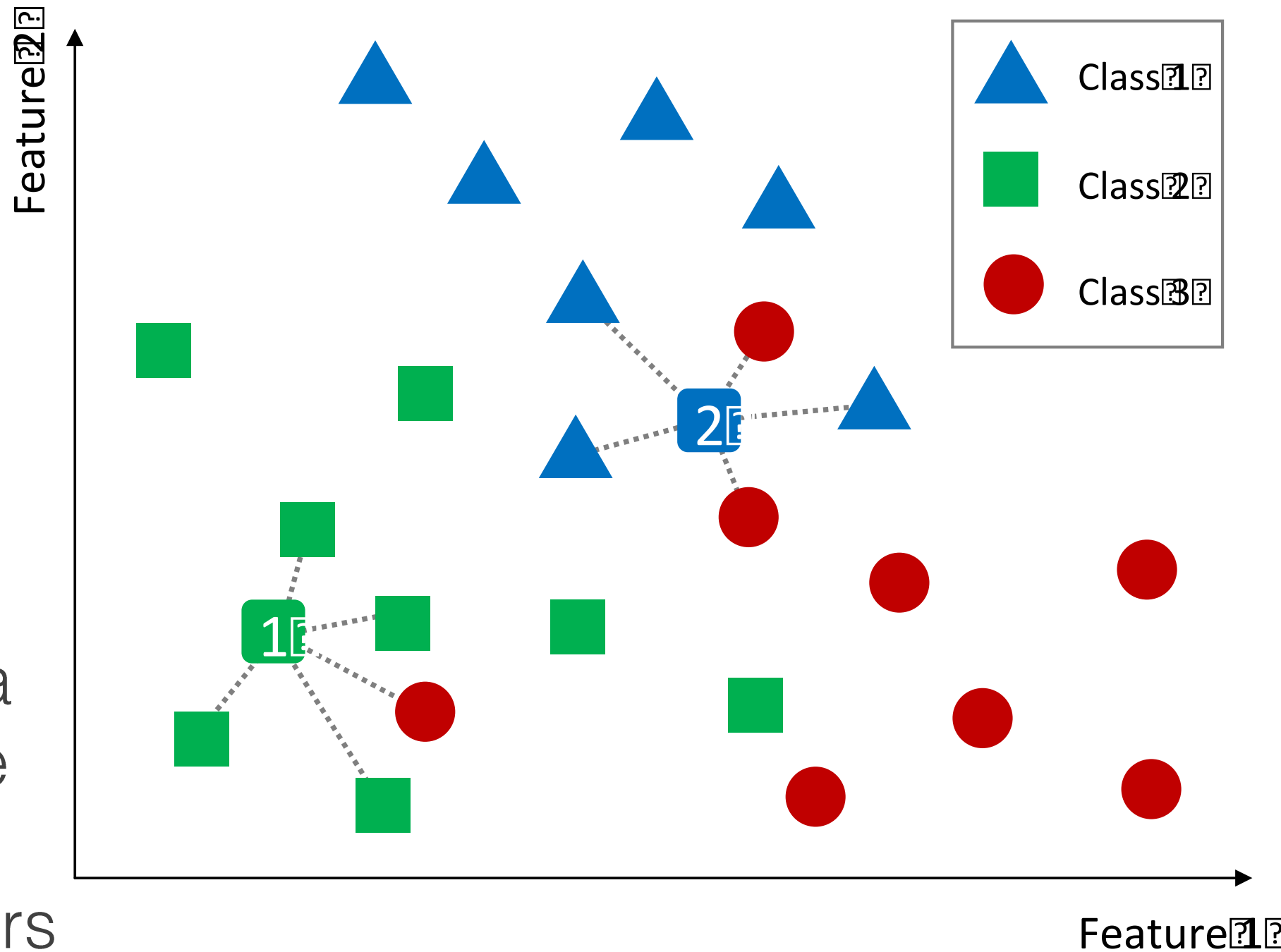
K Nearest Neighbor Classifier

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



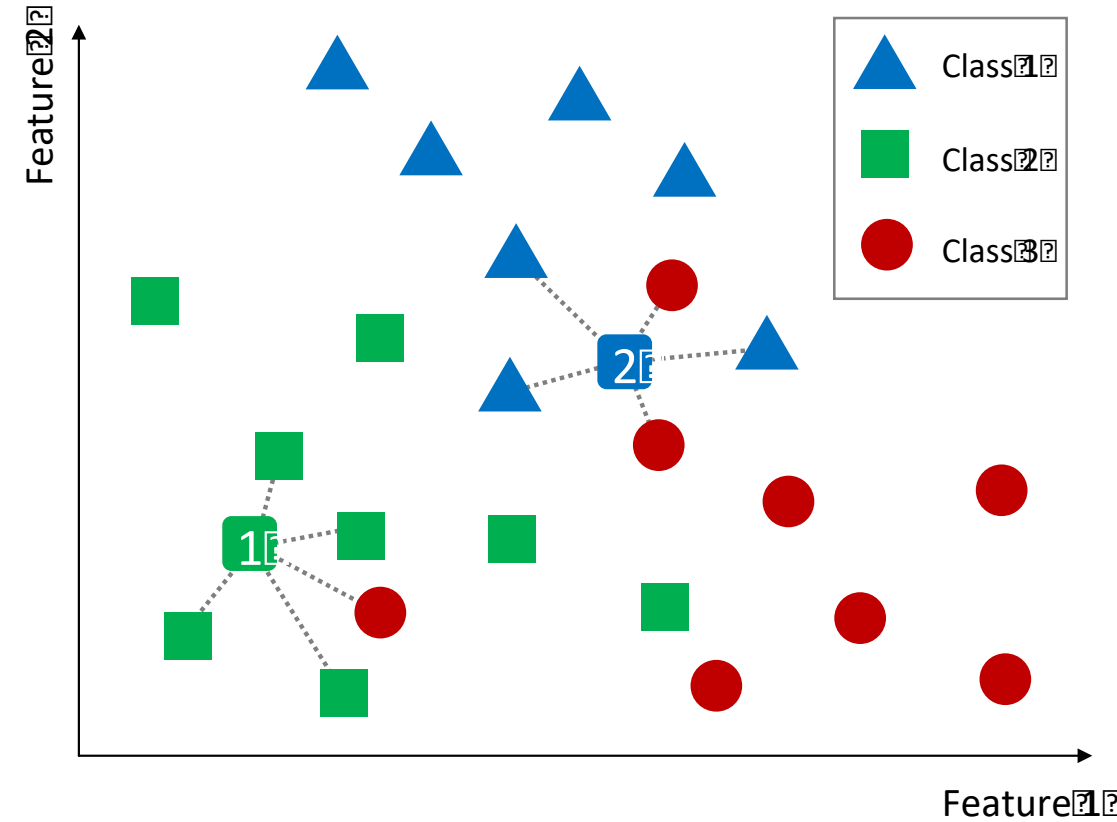
K Nearest Neighbor Classifier

Score vs Decision :

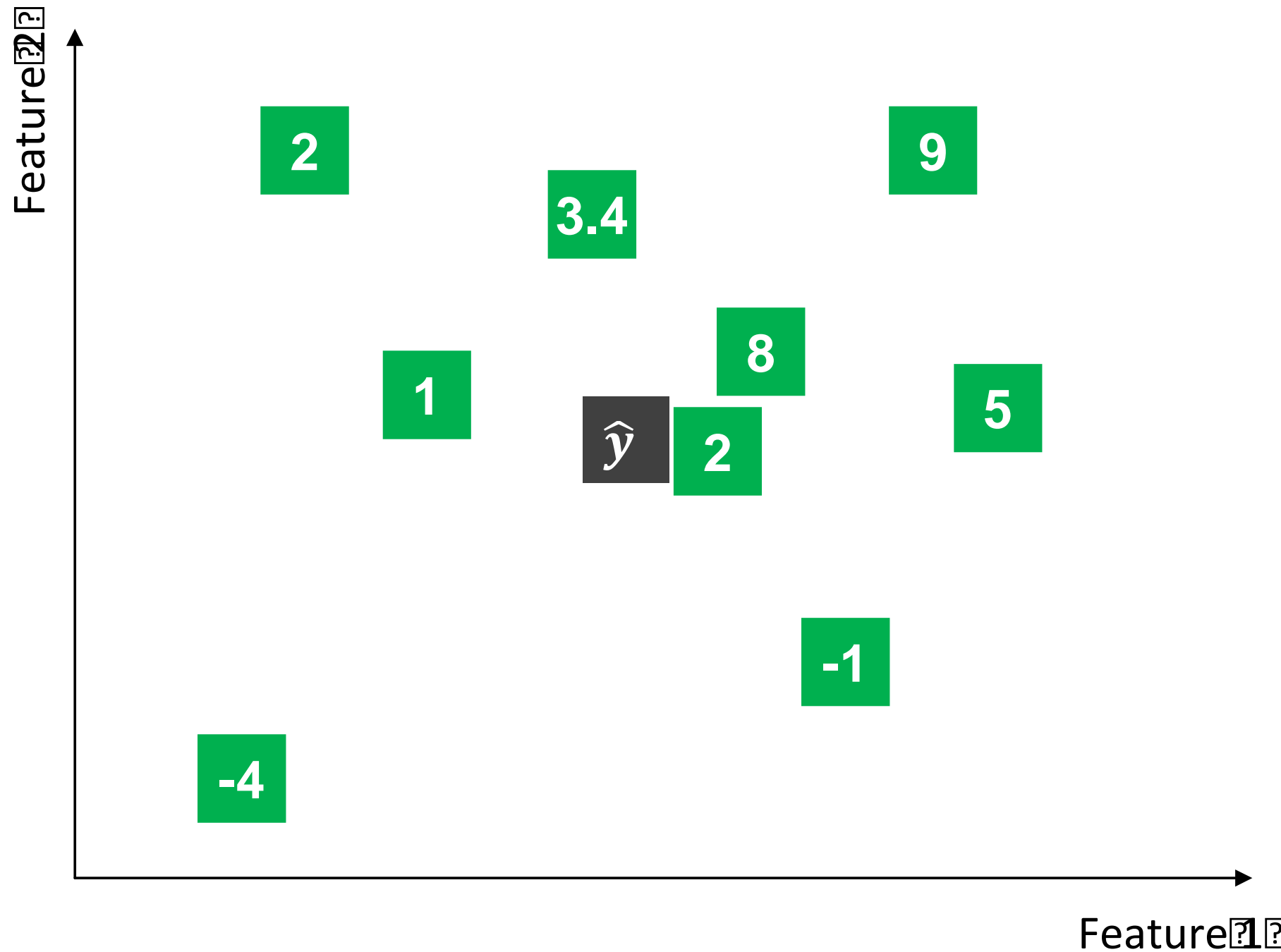
For 5-NN, the confidence score that a sample belongs to a class could be: $\{0, 1/5, 2/5, 3/5, 4/5, 1\}$

Decision Rule:

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

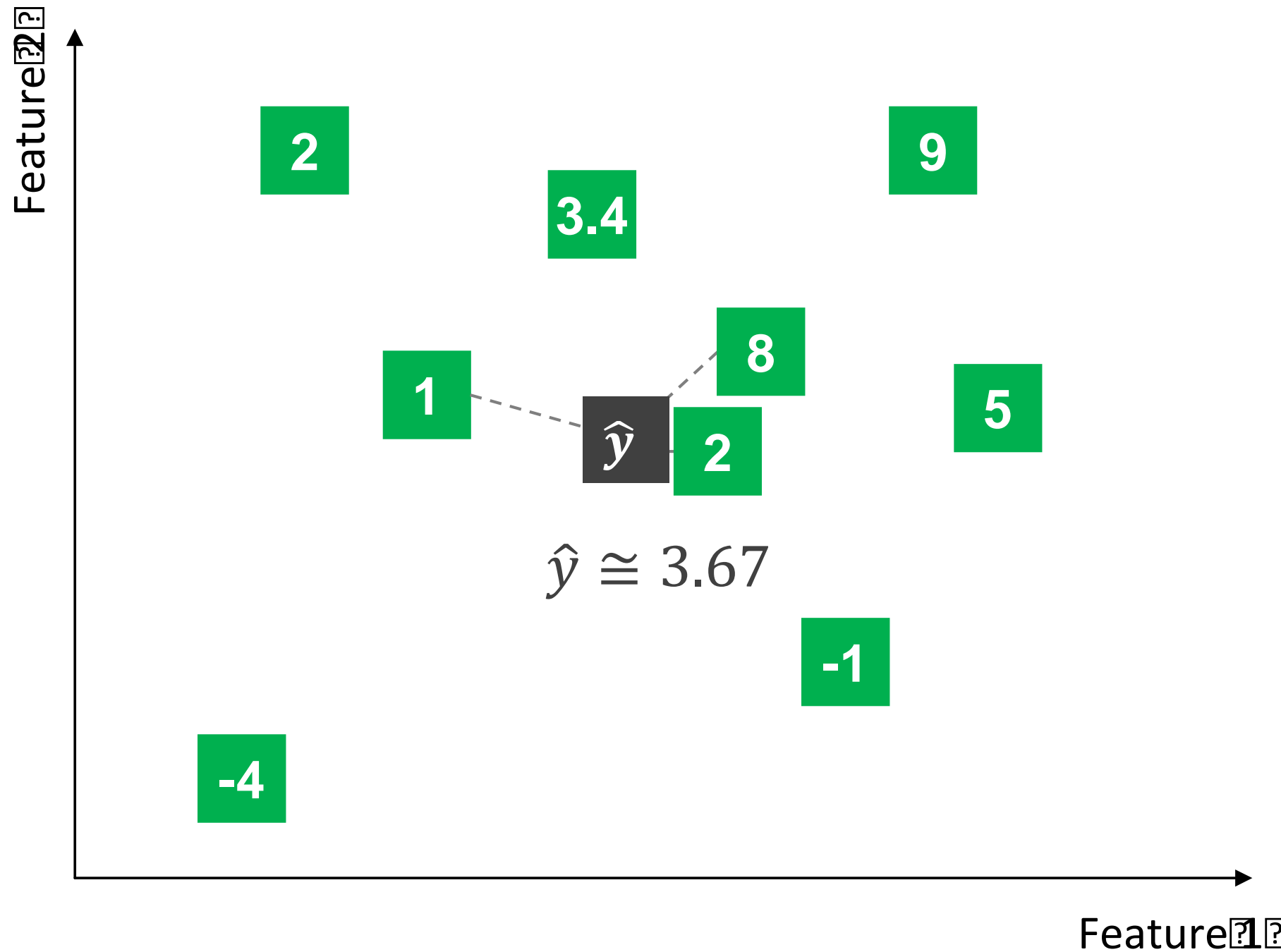


K Nearest Neighbor Regression



K Nearest Neighbor Regression

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



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

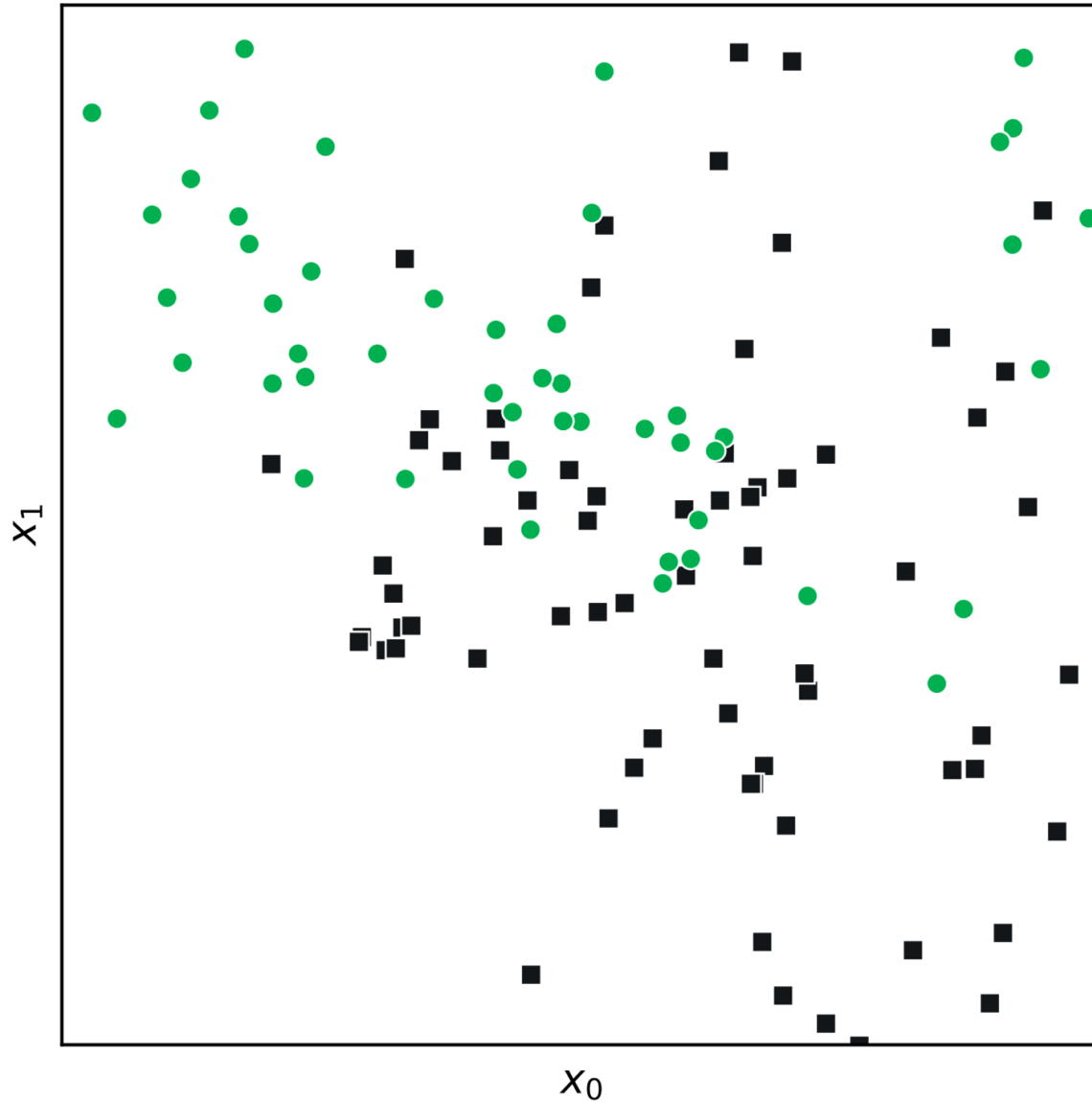
irreducible error inherent to the
problem

(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 = bias² + variance + noise

Classification feature space



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

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

Posterior

Likelihood Prior

Evidence

X Features

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

Bayes' Decision Rule:

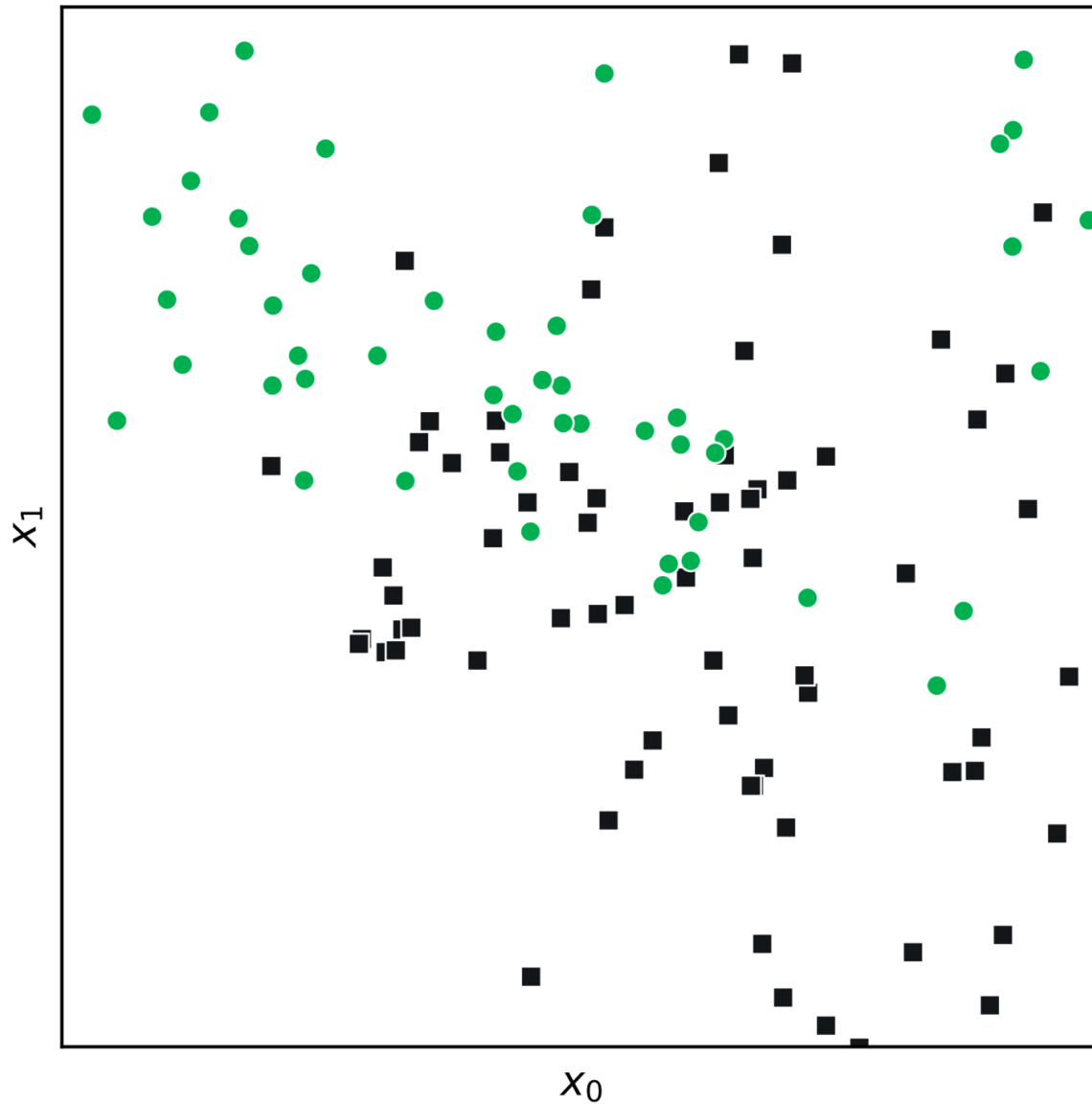
choose the most probable class given the data

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

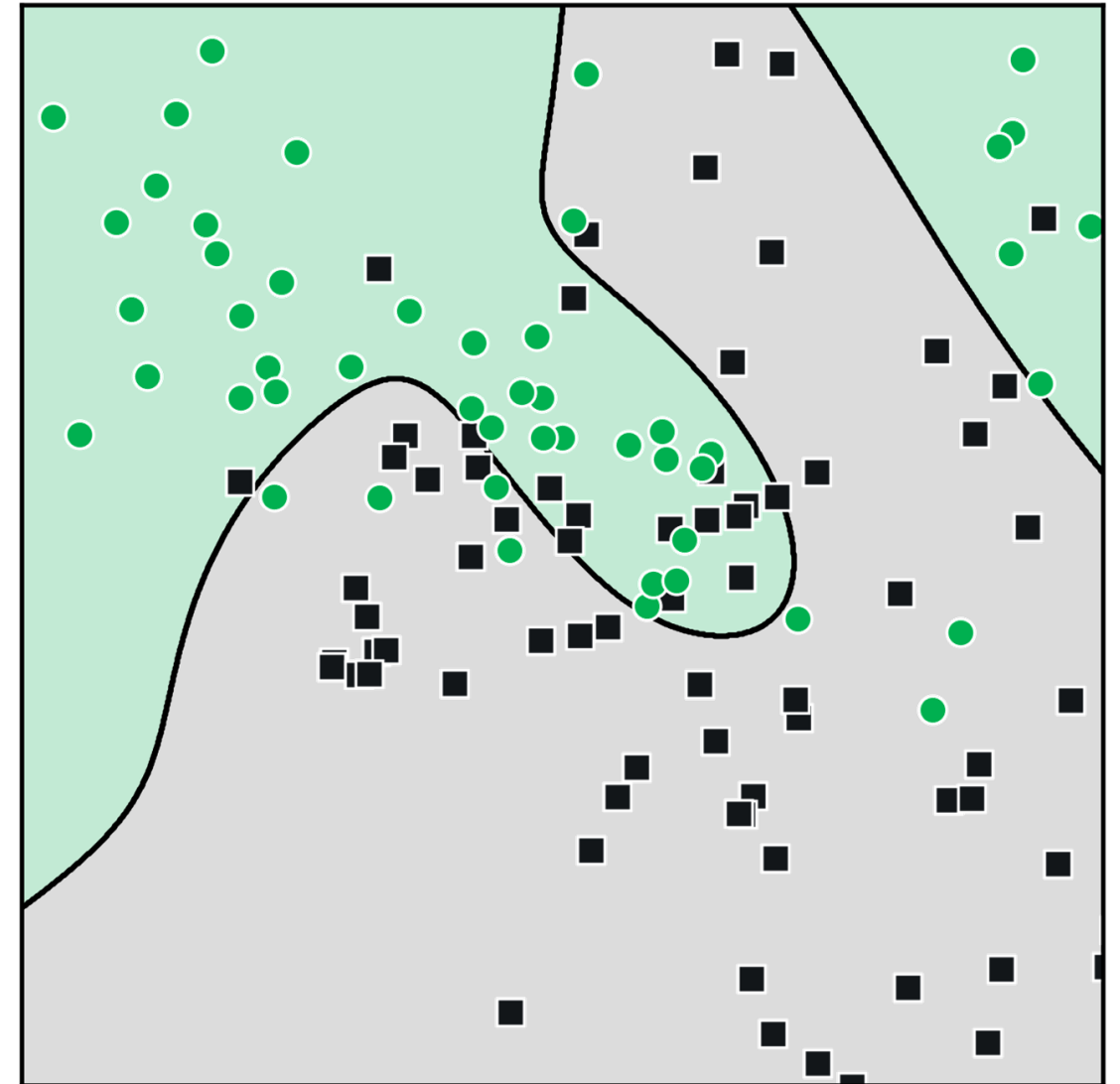
otherwise $\hat{y} = c_0$

- If the distributions are correct, this decision rule is **optimal**
- Rarely do we have enough information to use this in practice

Classification feature space

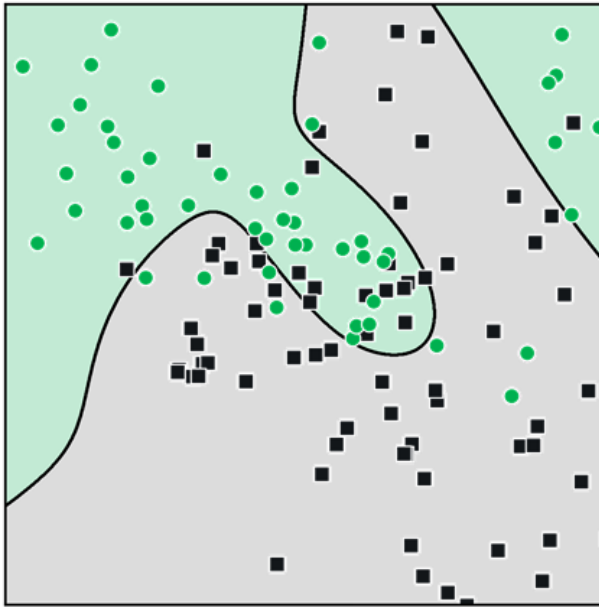


Bayes Decision Boundary

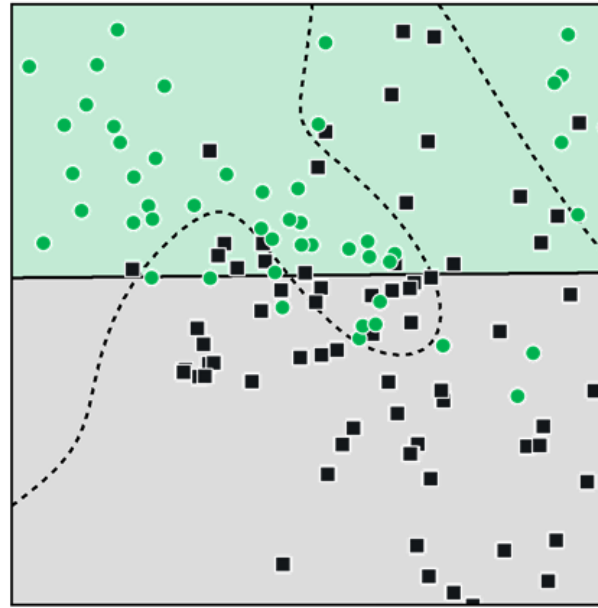


Decision Boundary Examples

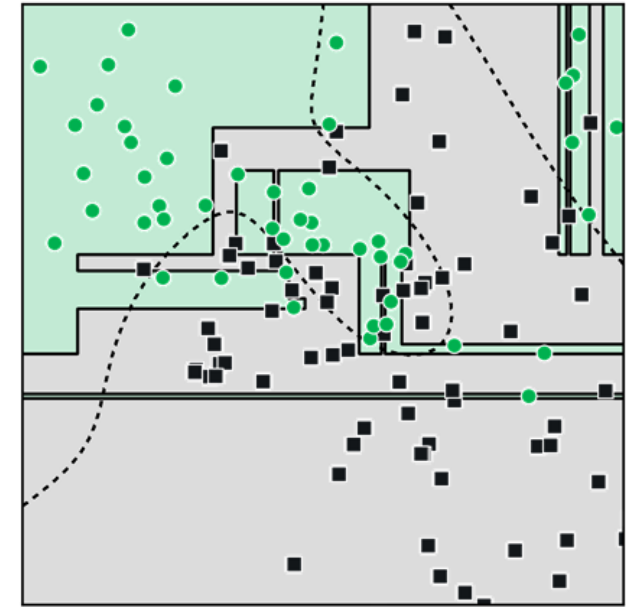
Bayes Decision Boundary



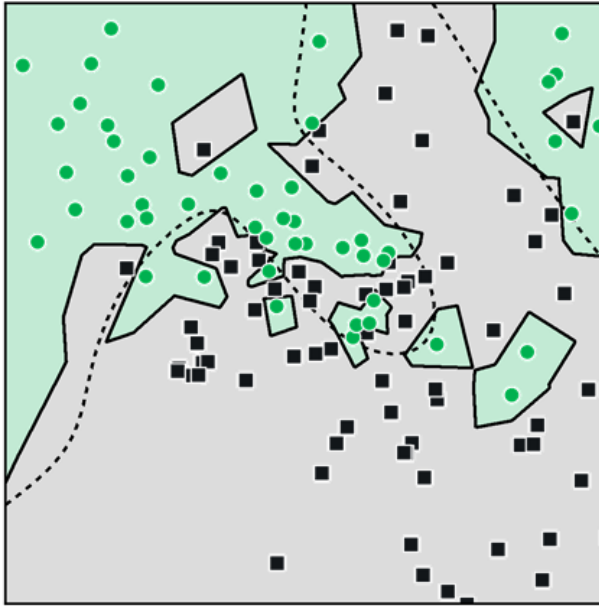
Linear Classifier



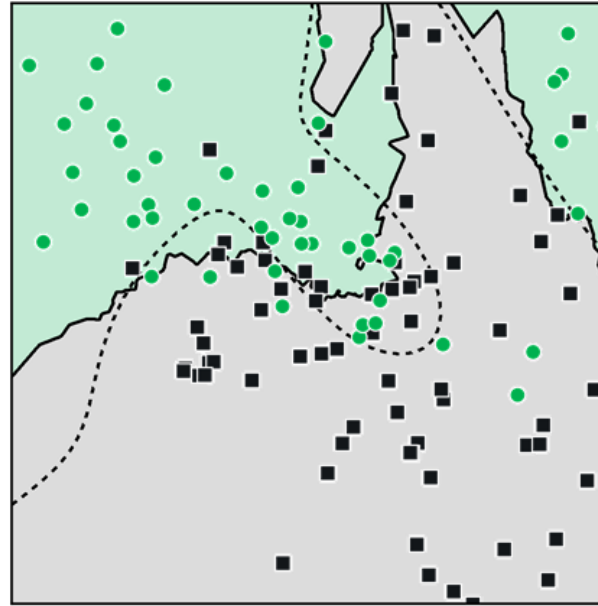
Decision Tree



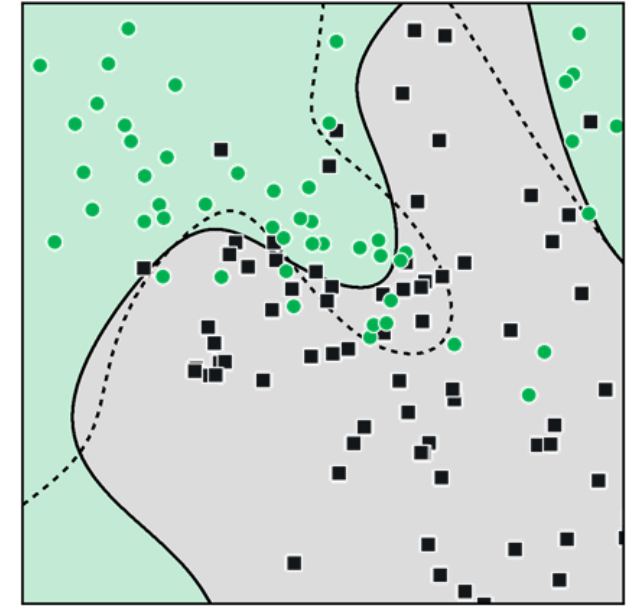
k=1 Nearest Neighbor



k=15 Nearest Neighbor



Support Vector Machine

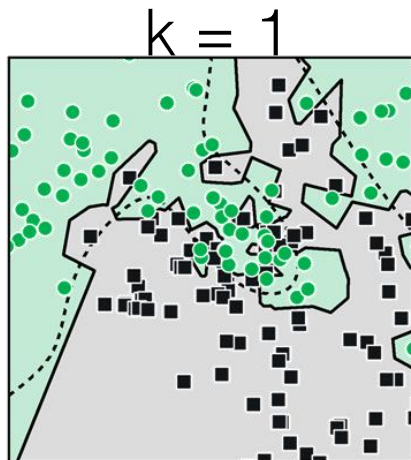


Bias Variance Tradeoff

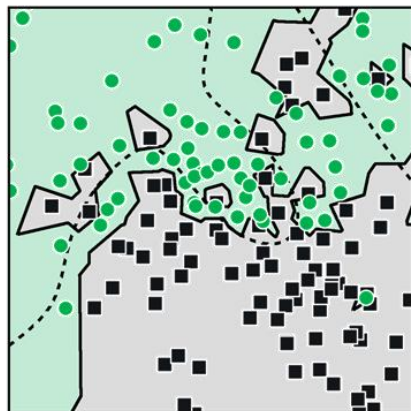
higher bias
underfit

higher variance
overfit

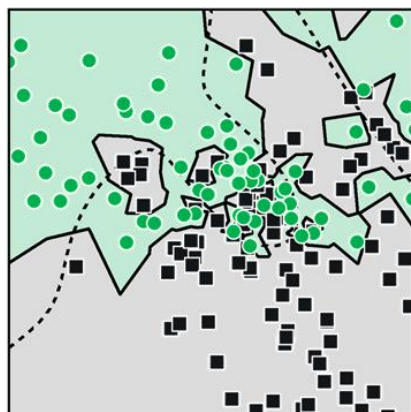
Sample 1



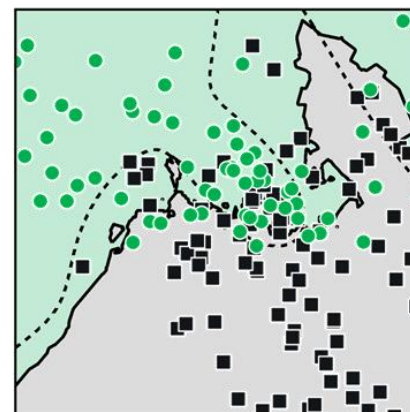
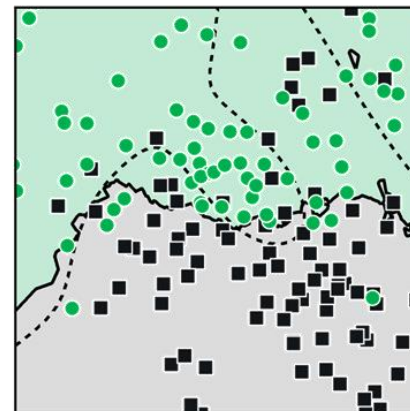
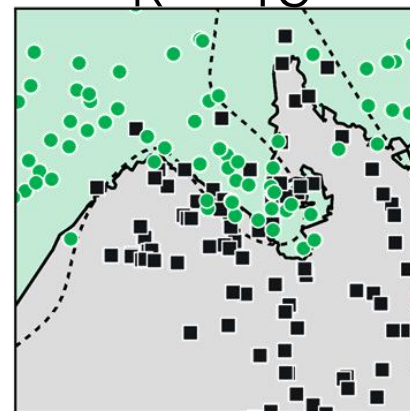
Sample 2



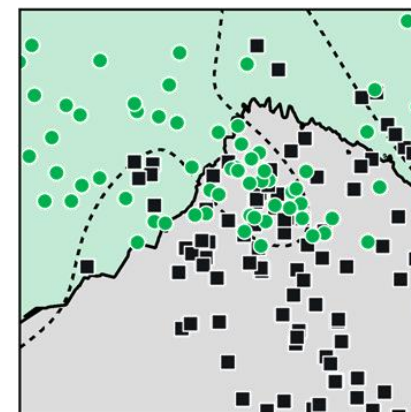
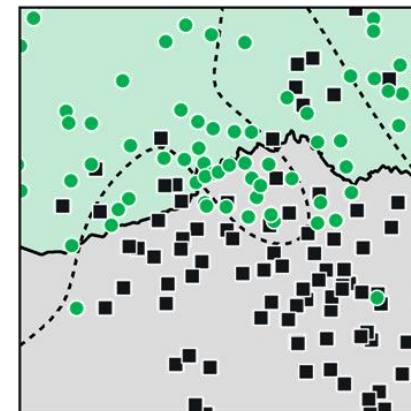
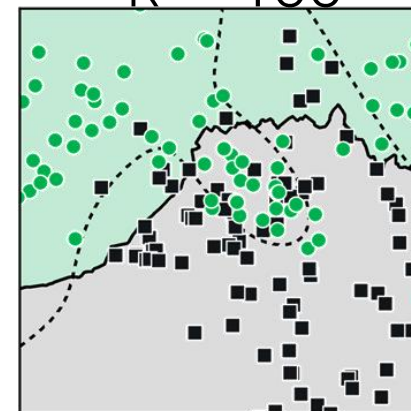
Sample 3



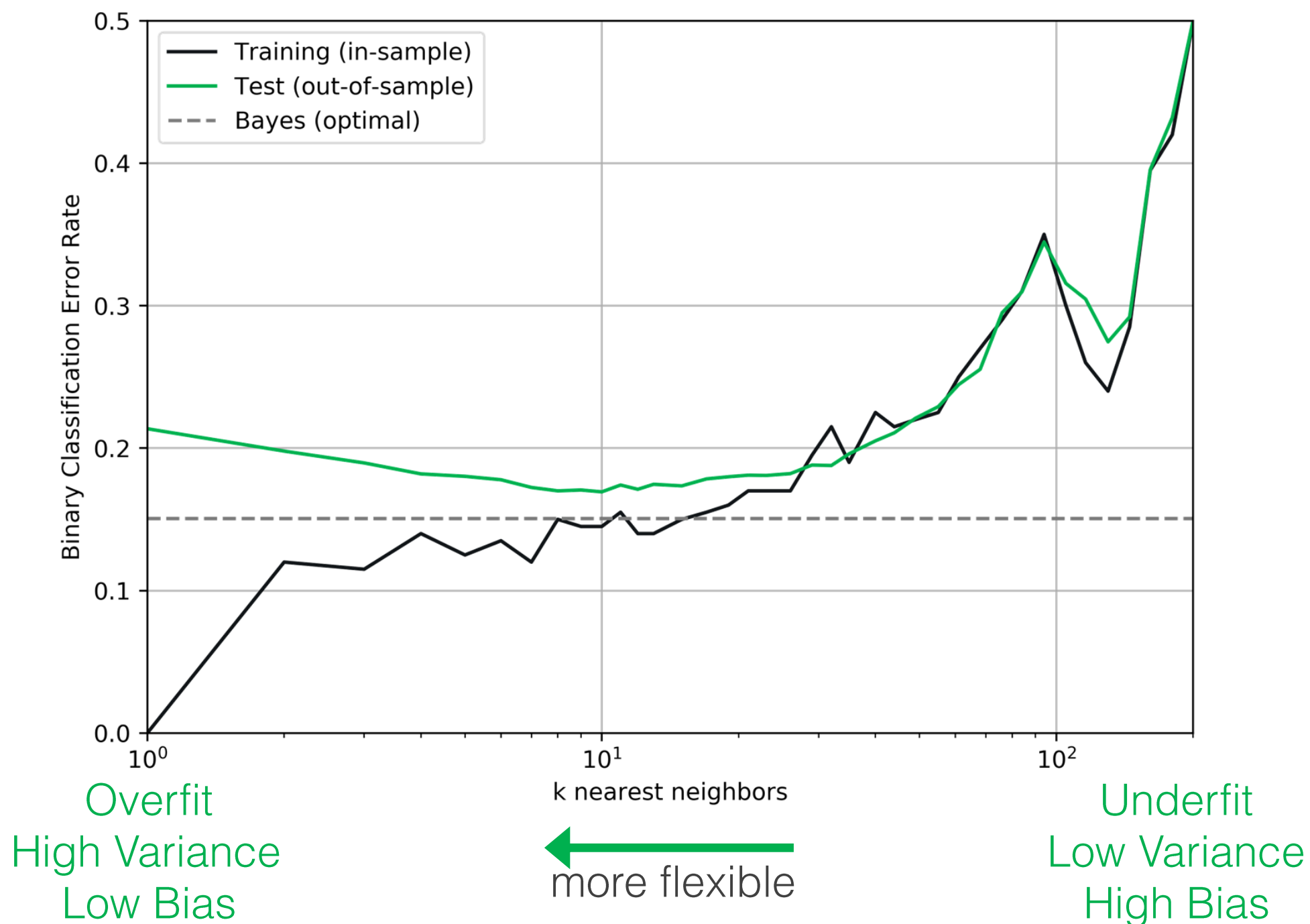
$k = 15$



$k = 100$



Bias Variance Tradeoff

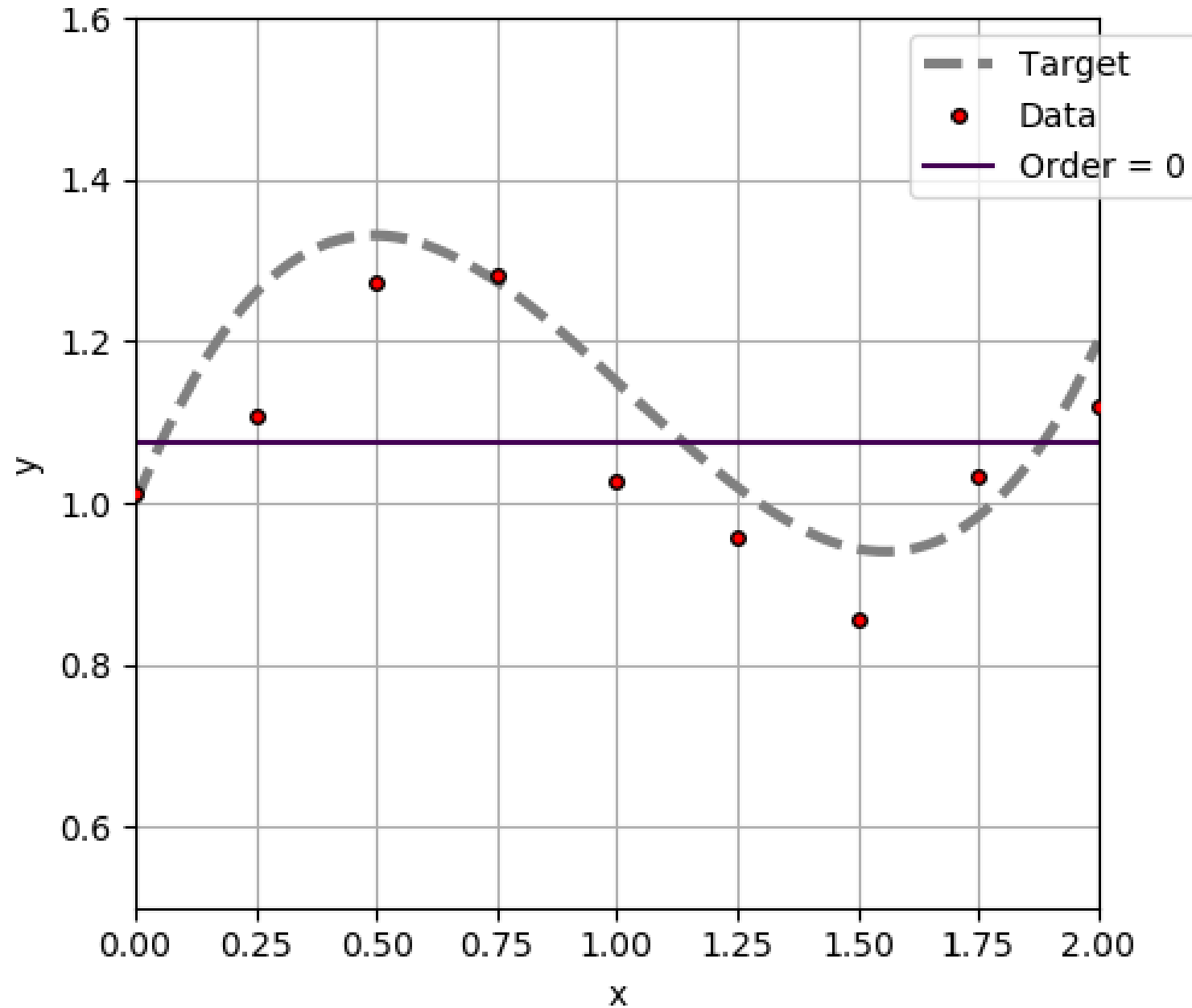


This tradeoff is equally challenging for regression

Linear Regression

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

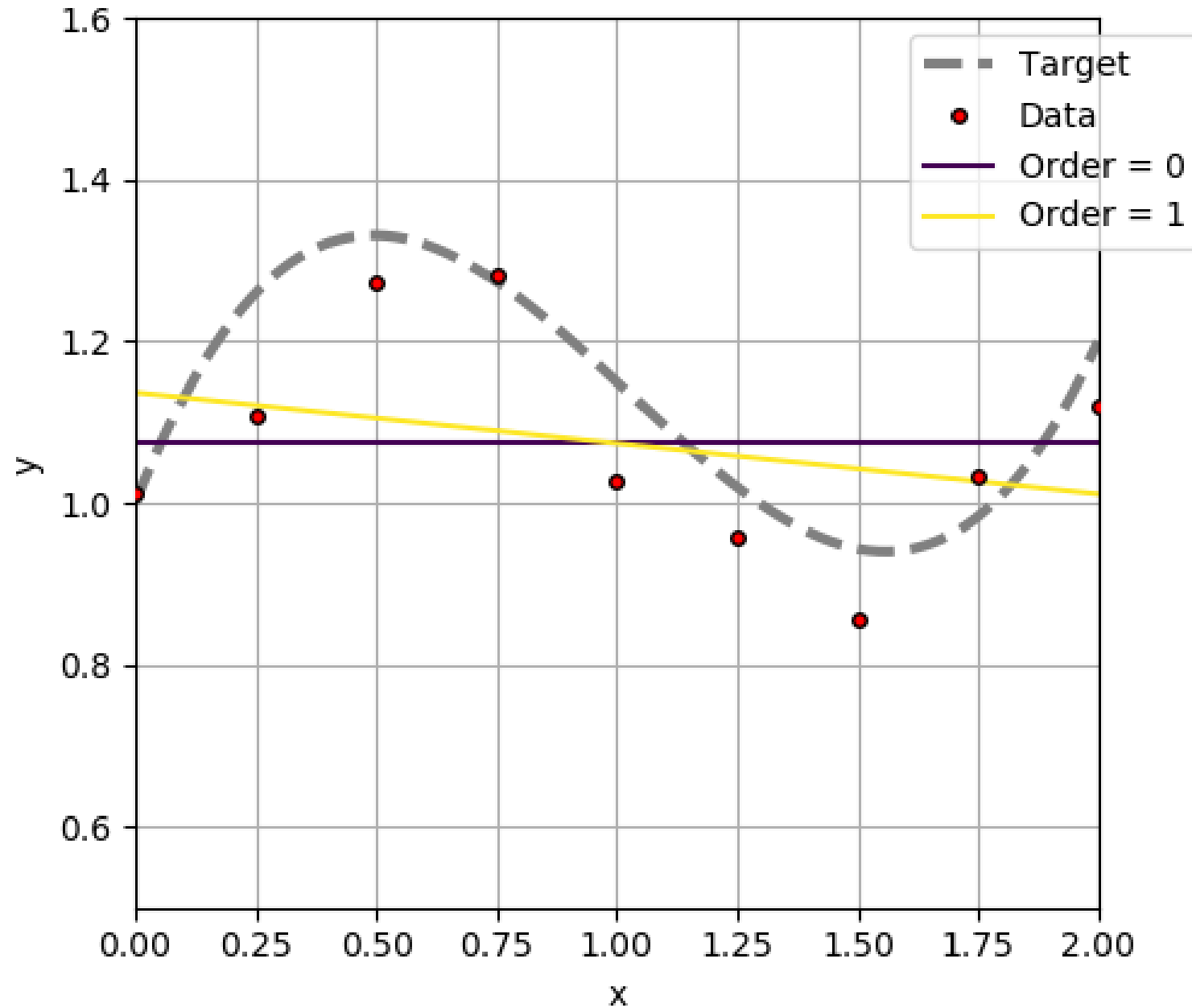
N is the model order



Linear Regression

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

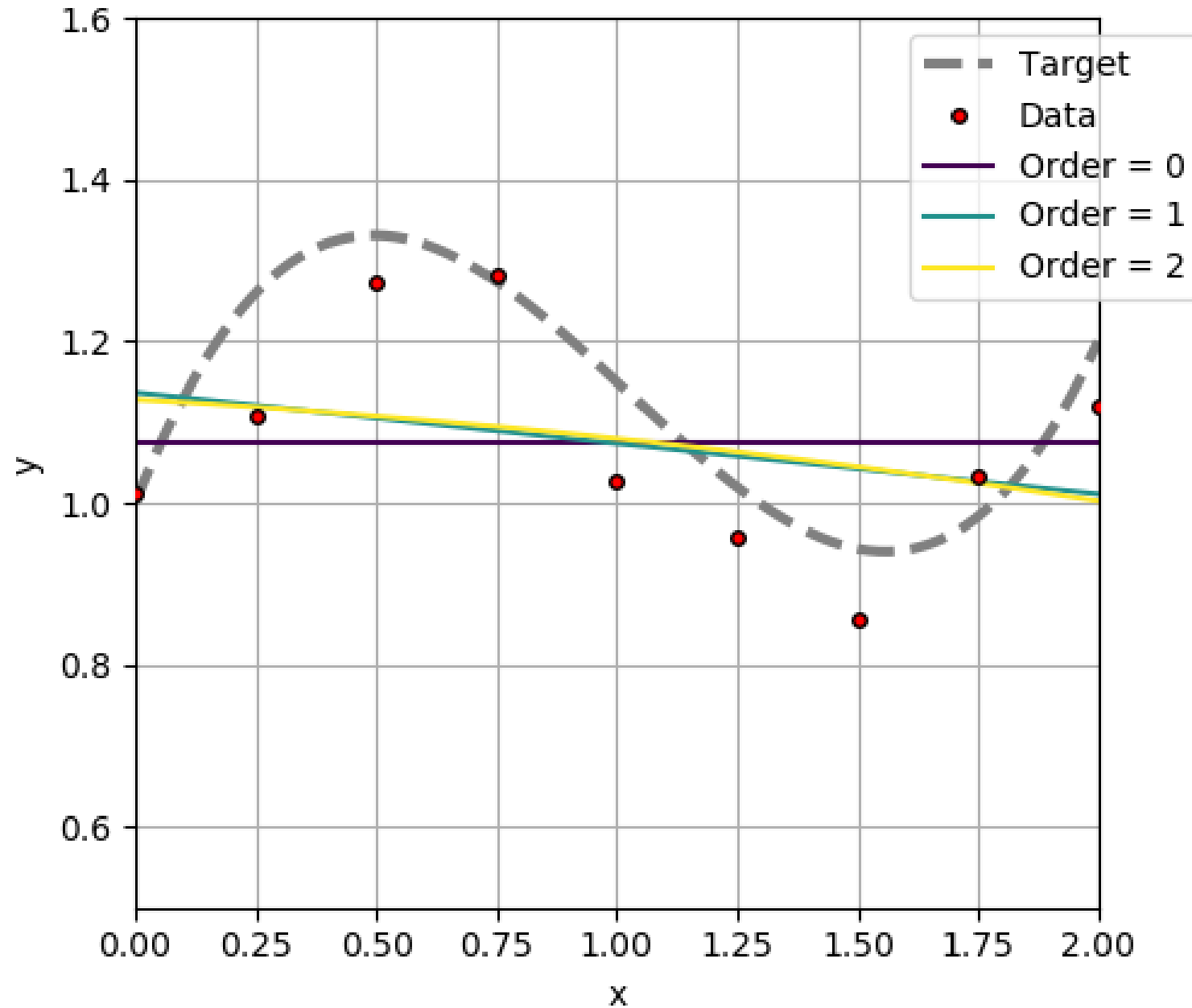
N is the model order



Linear Regression

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

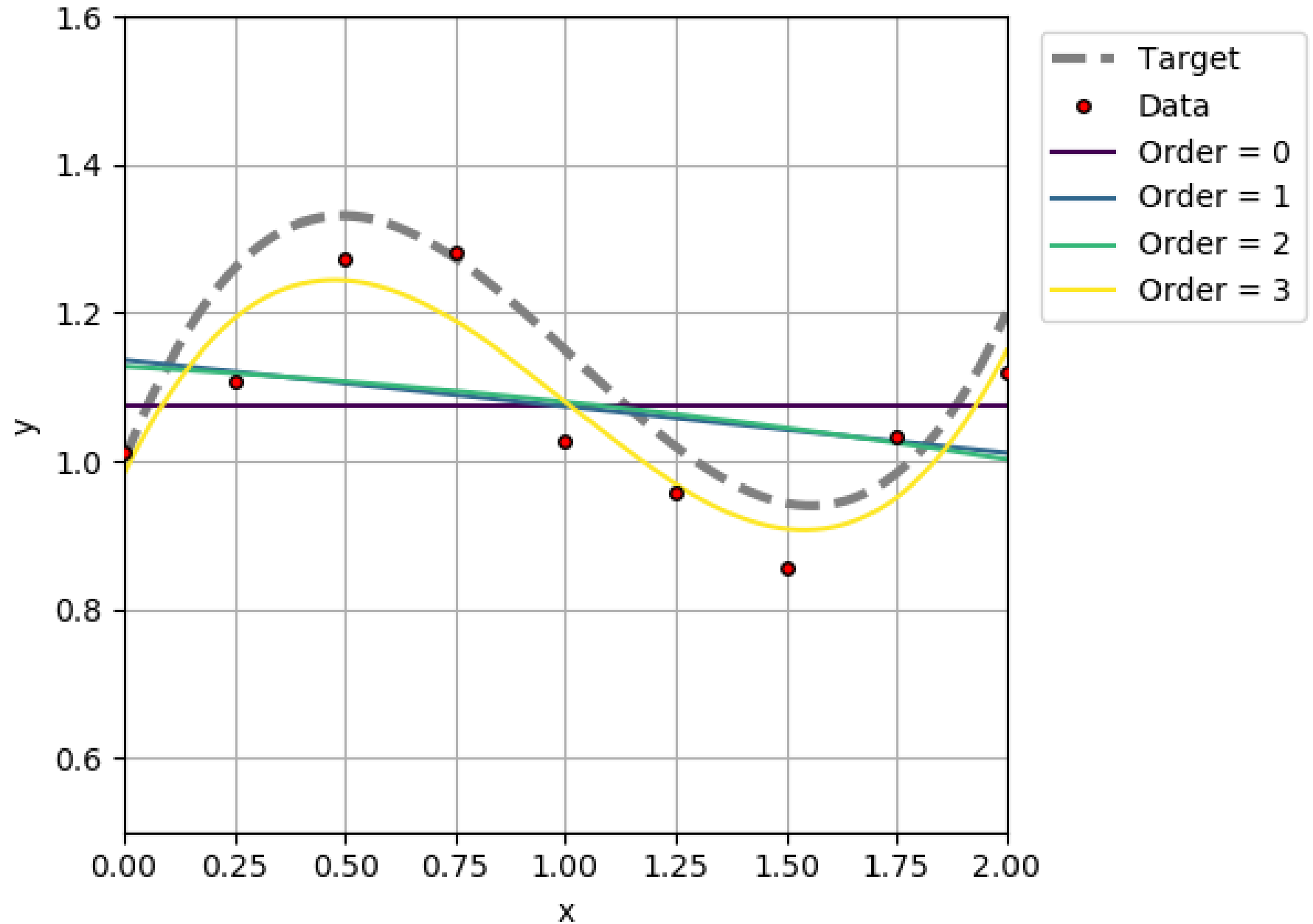
N is the model order



Linear Regression

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

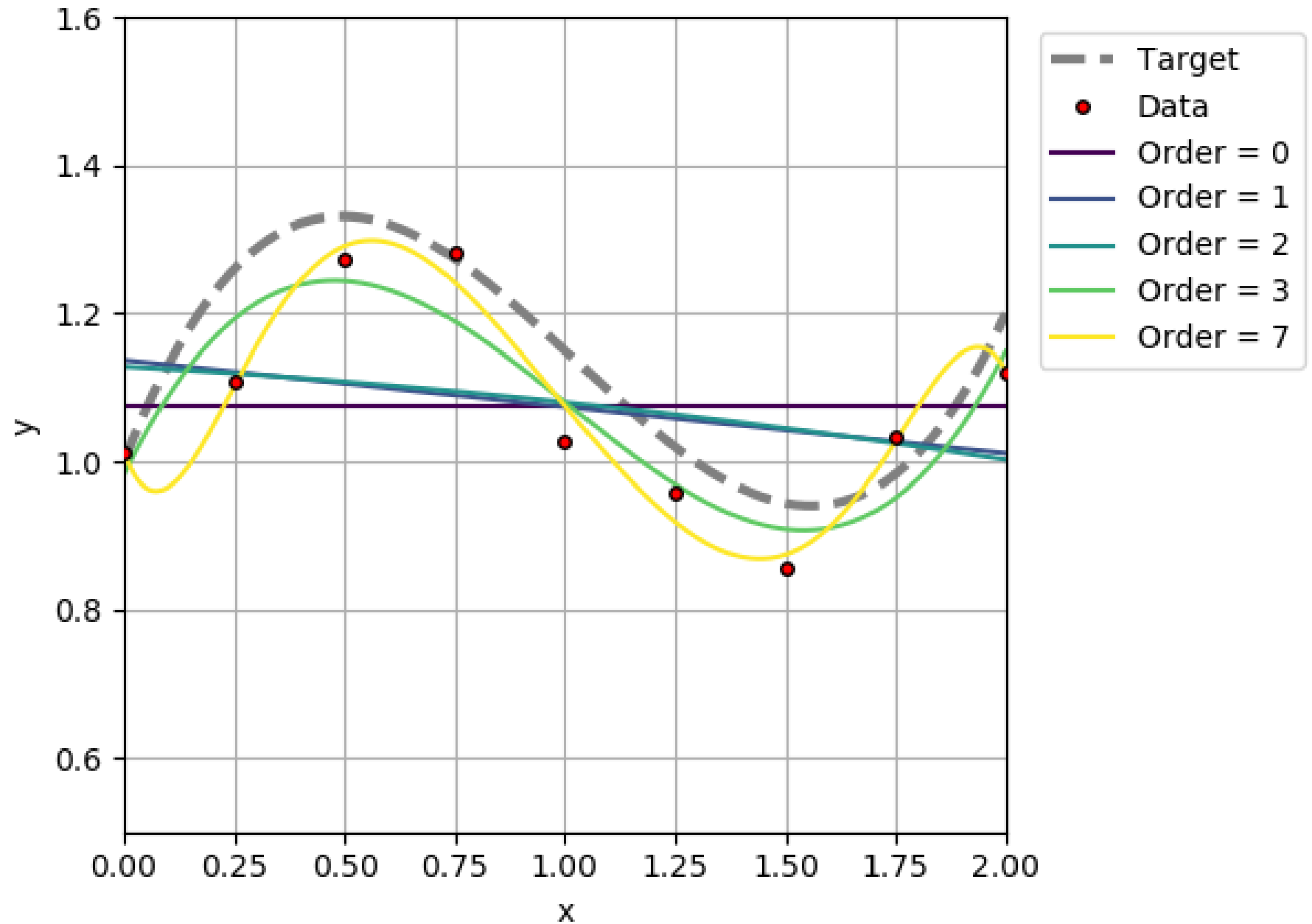
N is the model order



Linear Regression

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

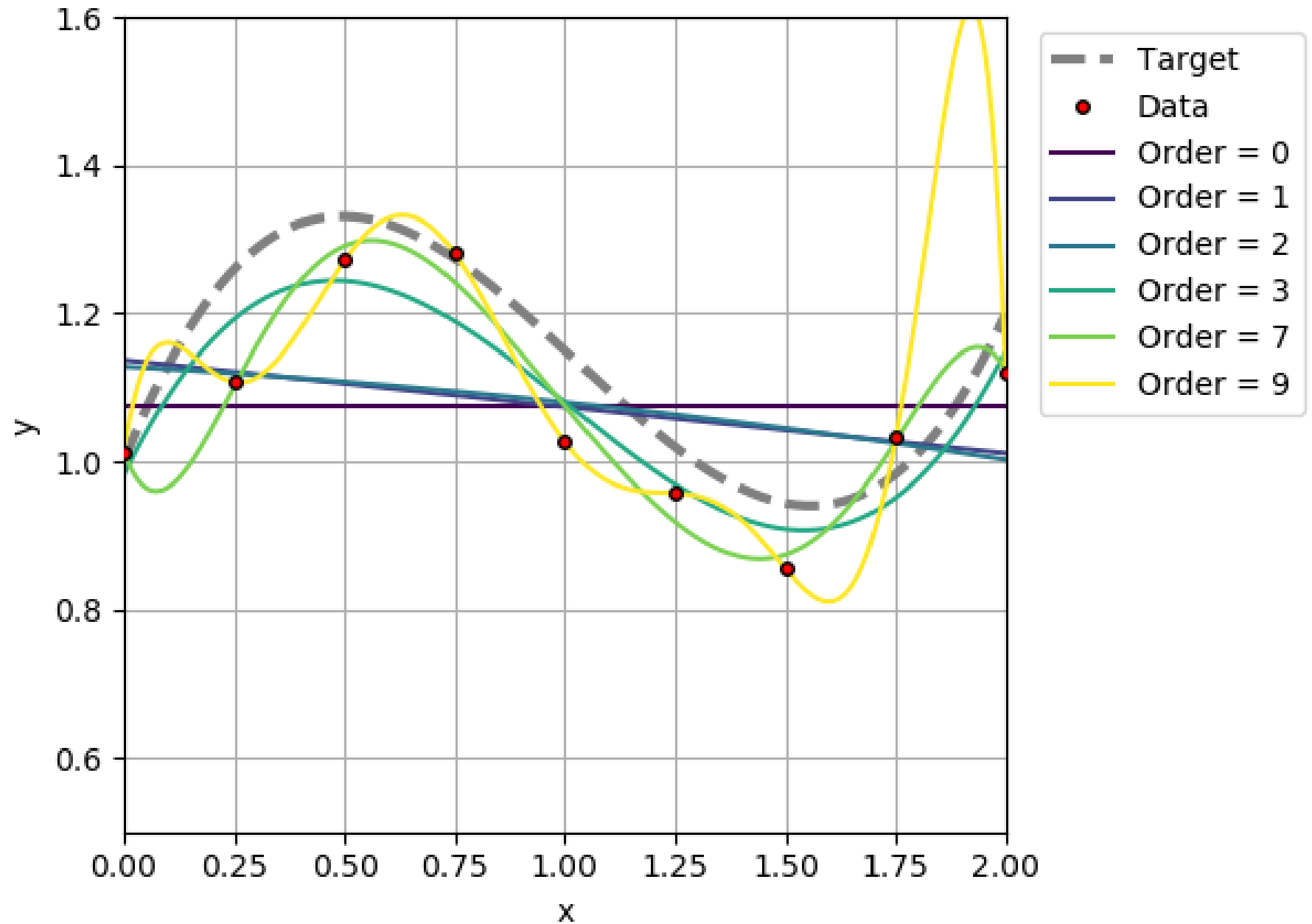
N is the model order



Linear Regression

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

N is the model order



Problem

Too much flexibility leads to **overfit**

Too little flexibility leads to **underfit**

Over/underfit **hurts generalization** performance

Solutions

1. Add more data for training
2. Constrain model flexibility through **regularization**