

Lecture 5: Machine learning

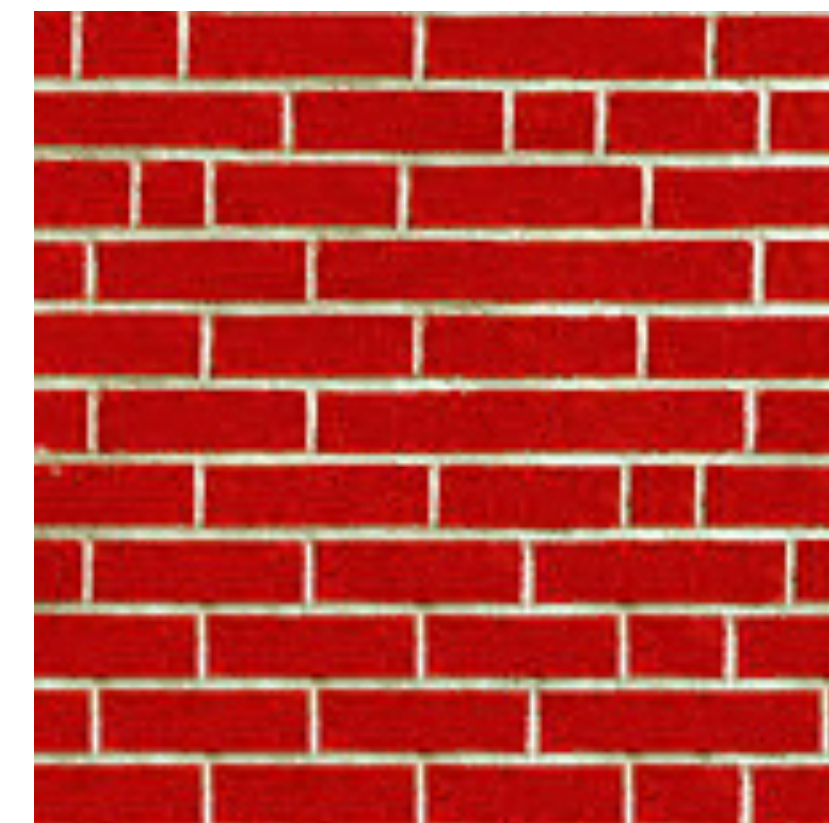
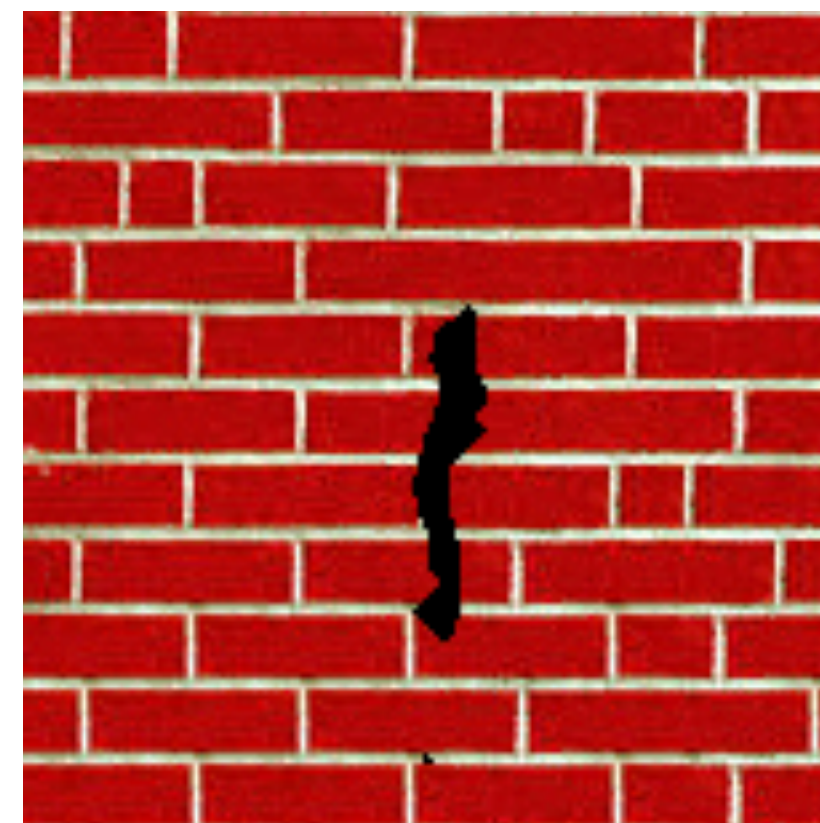
Announcements

- PS1 due tonight
- PS2 out tonight (by midnight).
- Due 1 week from today!
- Szeliski 5.1
- Optional: Goodfellow Deep Learning “ML Basics

Announcements

- Late days → Late hours
- You have 5×24 late hours total
- We no longer “round up” (e.g., turning in 1 hour late doesn’t take away a full day).
- Note that Piazza is a discussion platform, not just Q&A with course staff.
- Thank you for those who have been answering questions!
- Course staff may delay answering Piazza questions and may send you to office hours for complex questions.

Last week: inpainting (hole filling)



Simple inpainting [Efros & Leung 1999]



[Hays and Efros “Scene Completion Using Millions of Photographs”, 2007.]





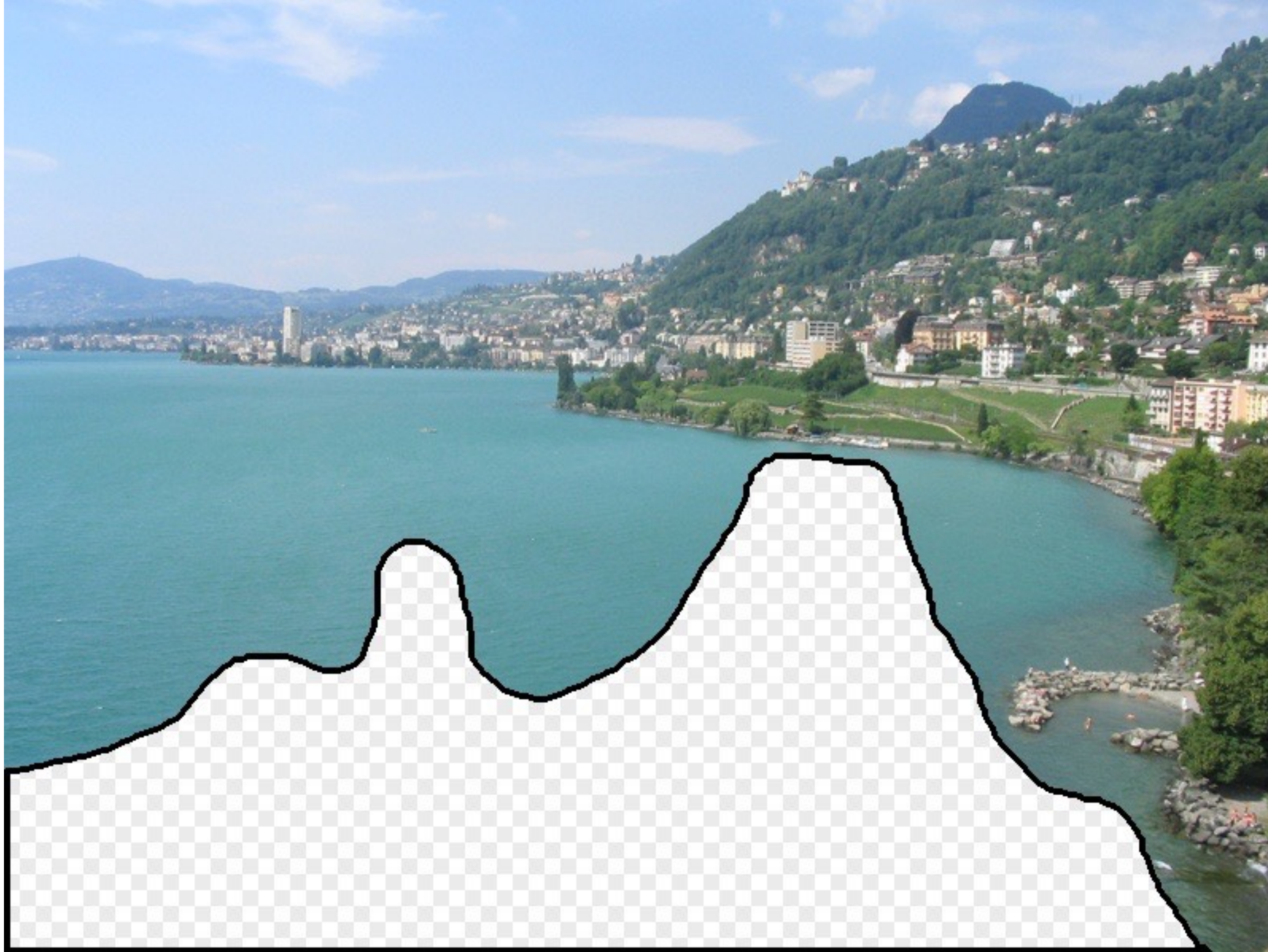
Simple inpainting [Efros & Leung 1999]

How would a human solve this?



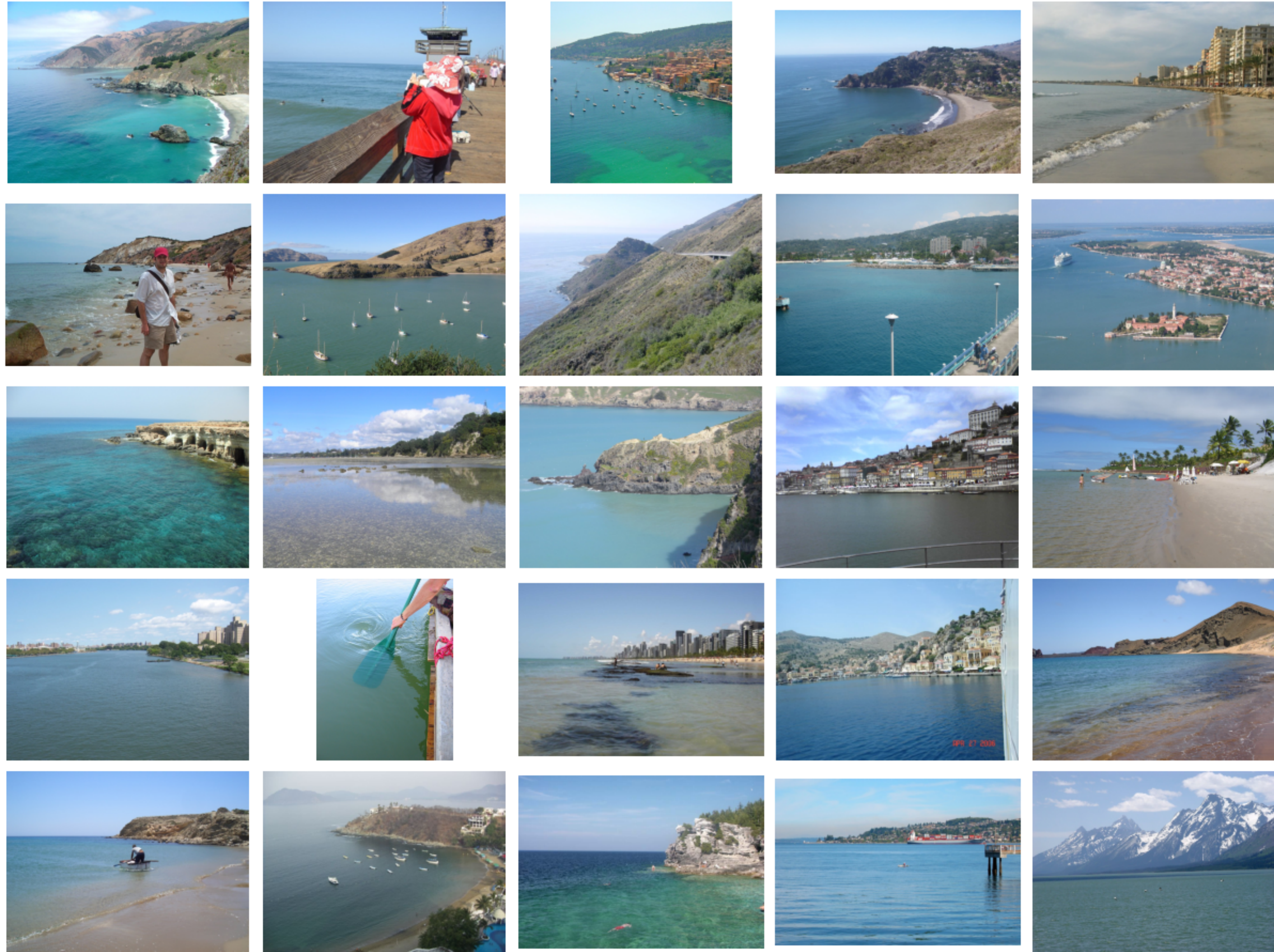
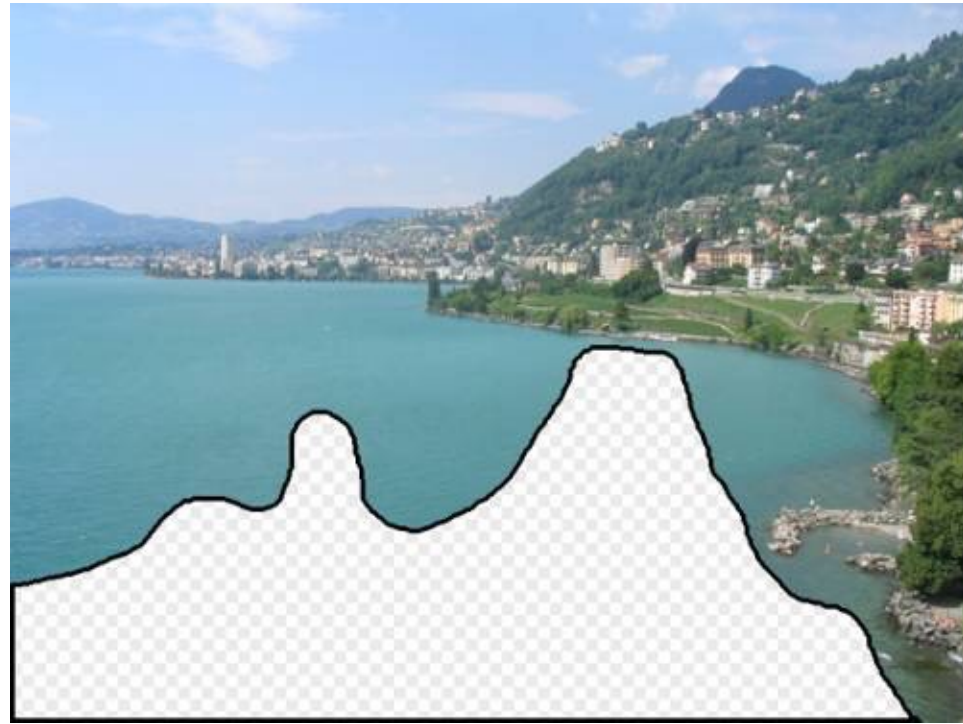


[Hays and Efros. Scene Completion Using Millions of Photographs. SIGGRAPH 2007.]

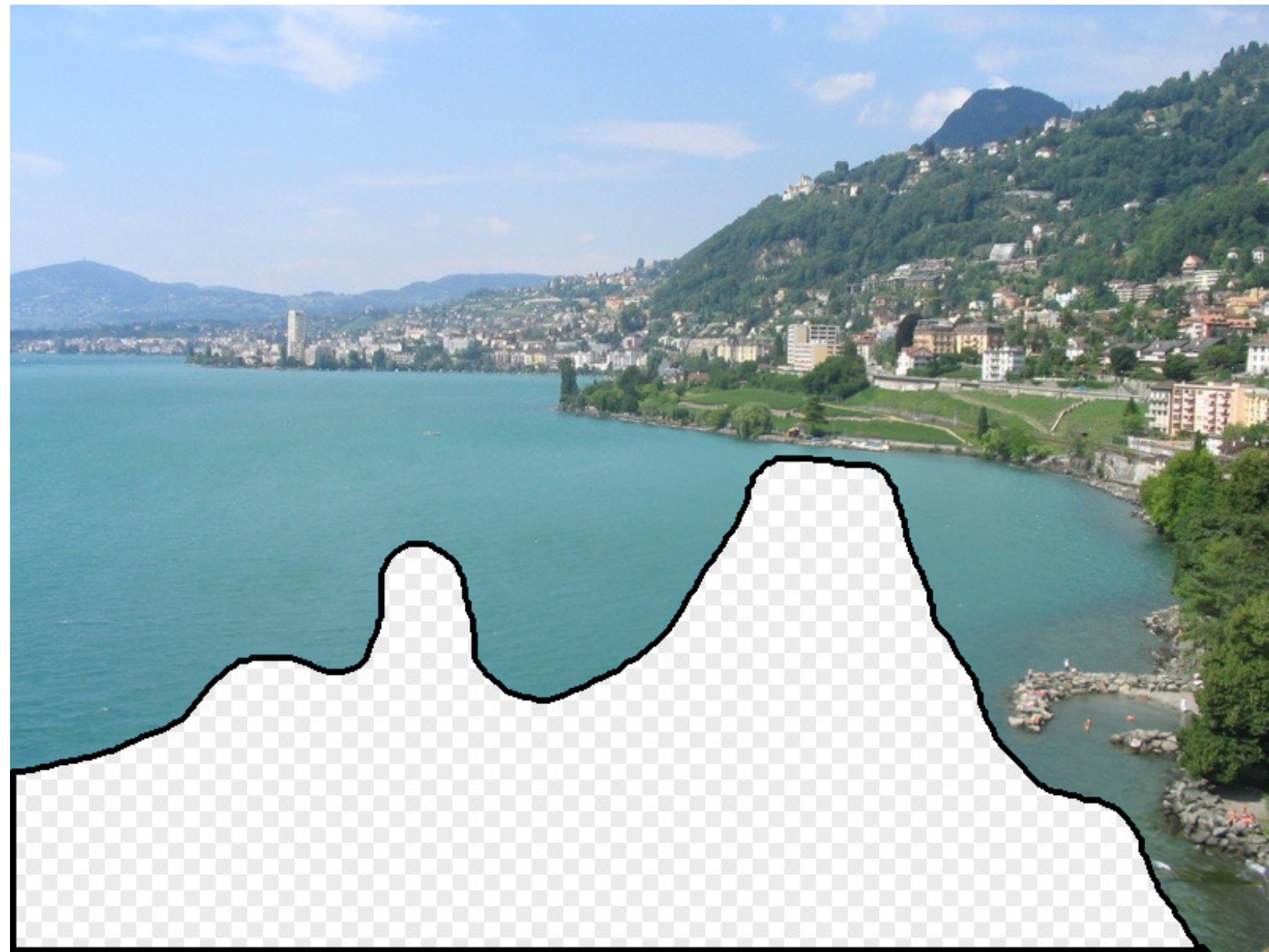


2 Million Flickr Images





... 200 total





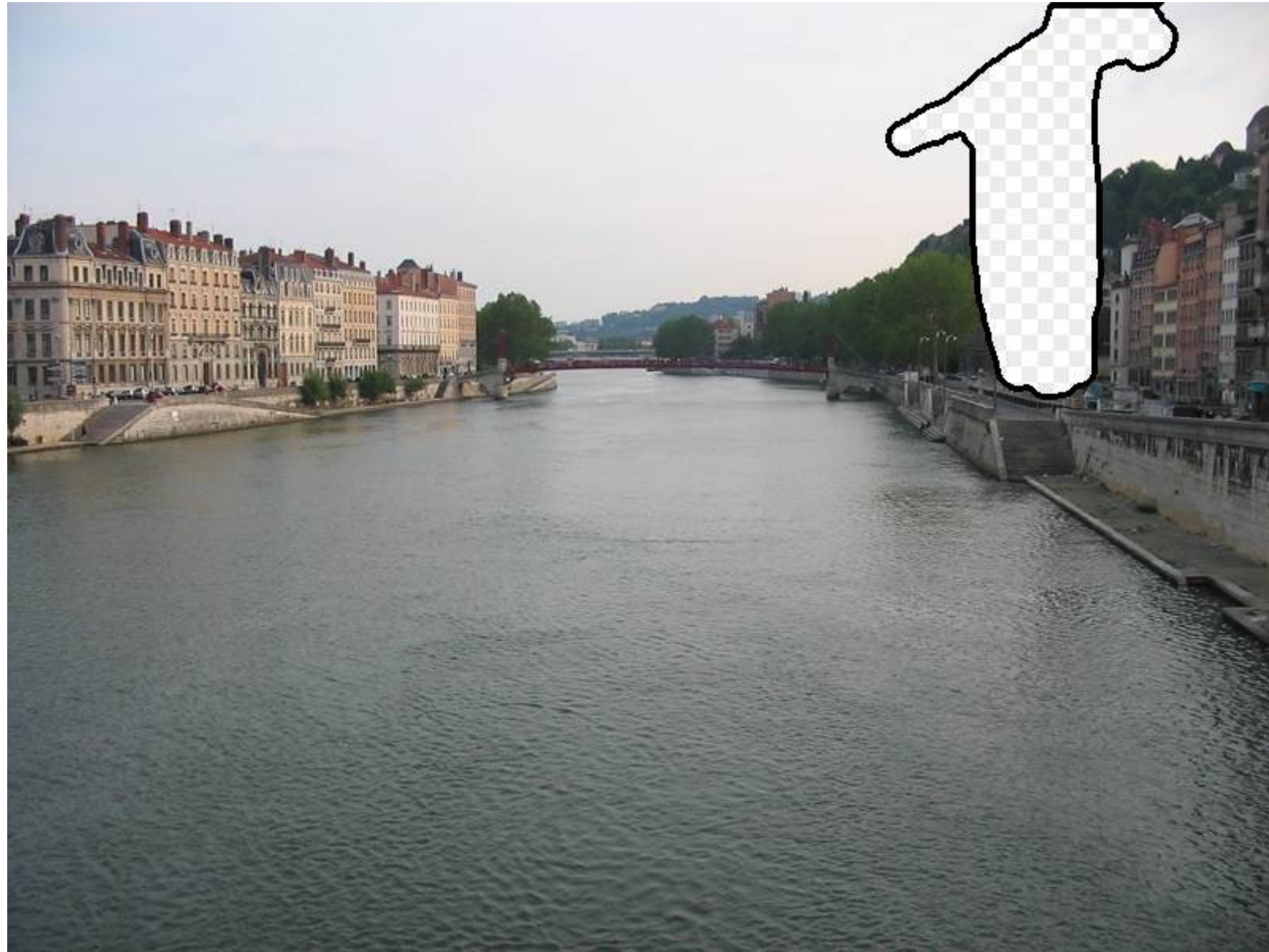




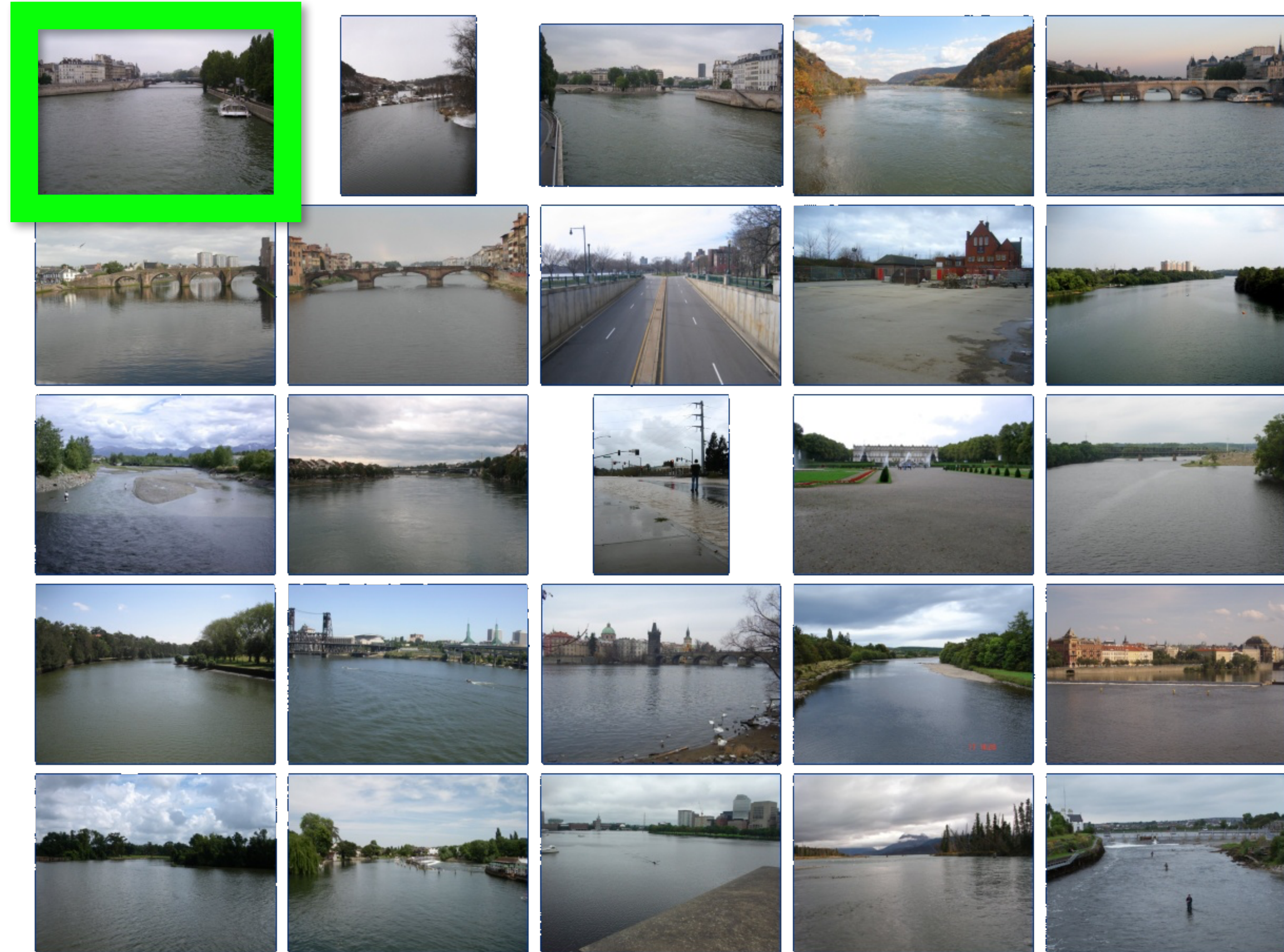
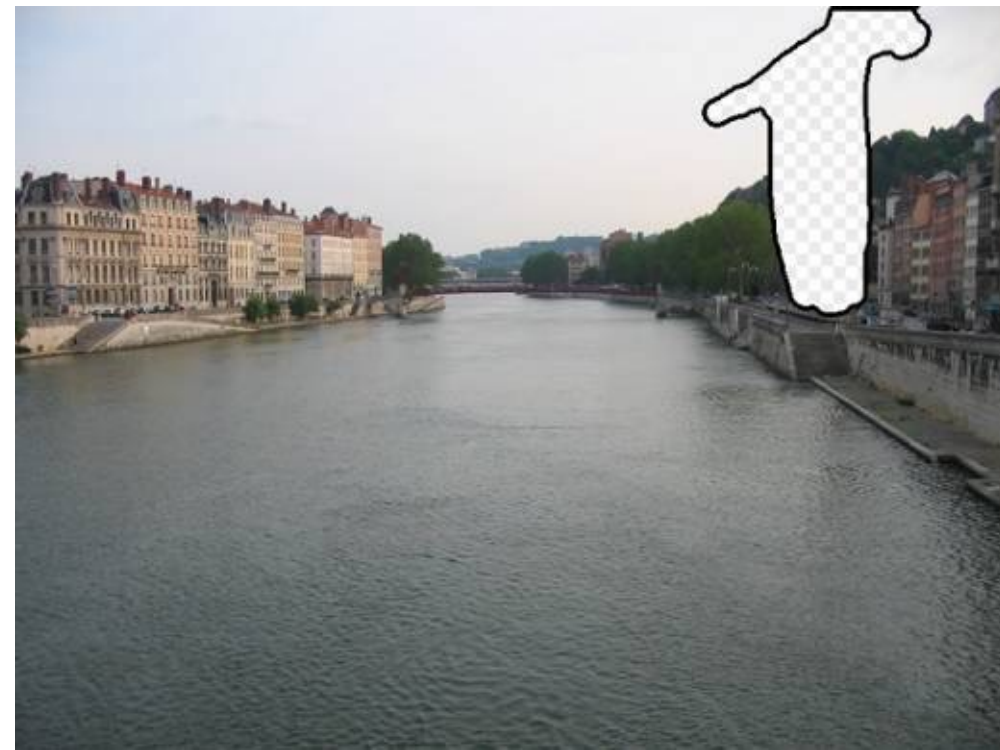












... 200 scene matches

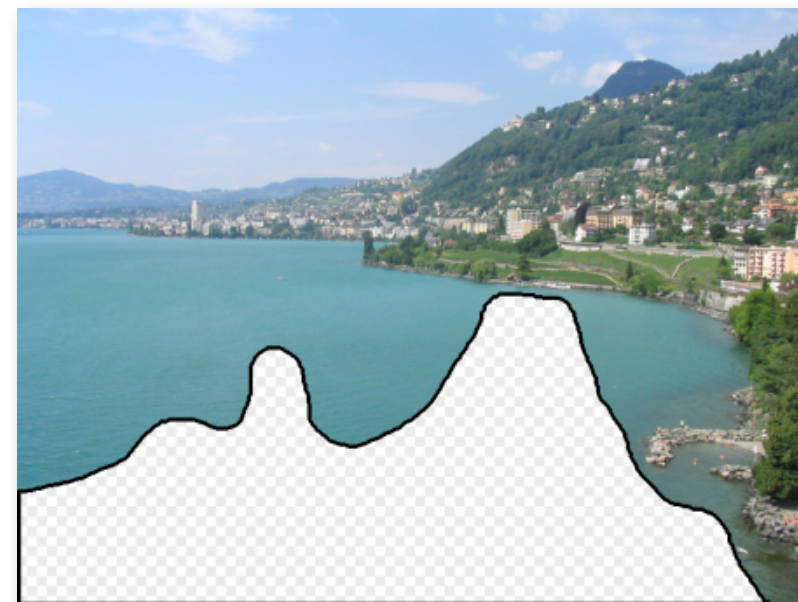


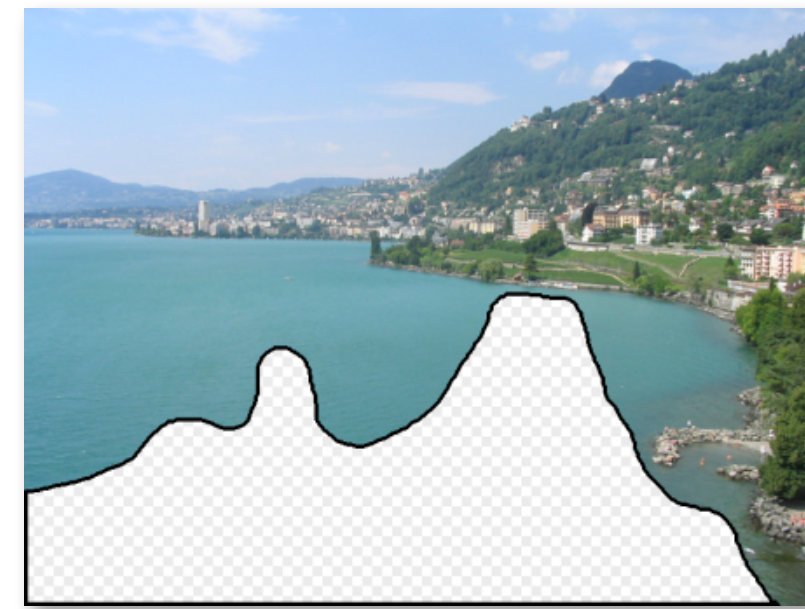




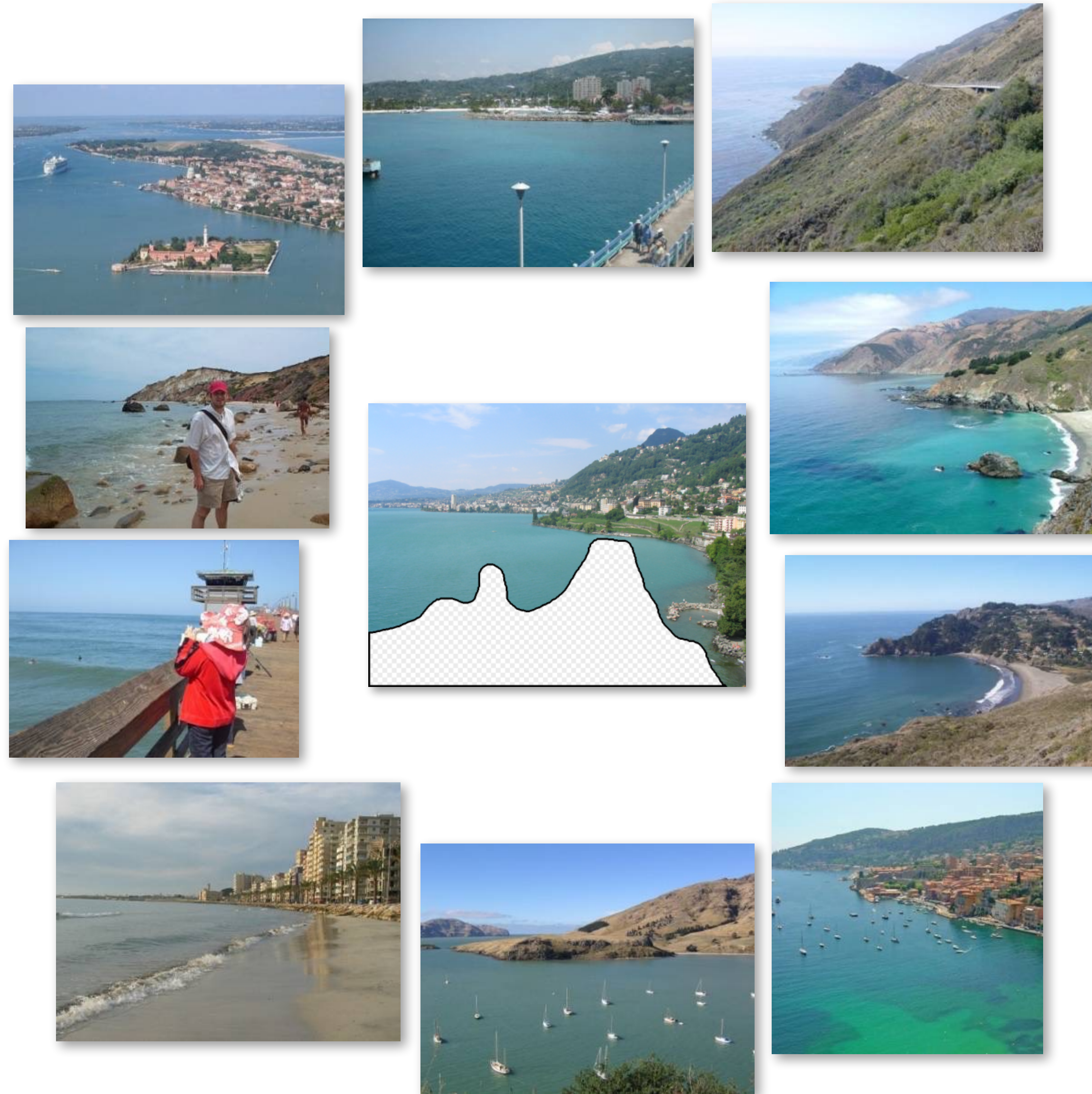


Why does this work?





Nearest neighbors from a collection of 20 thousand images



Nearest neighbors from a collection of 2 million images

“Unreasonable Effectiveness of Data”

[Halevy, Norvig, Pereira 2009]

Parts of our world can be explained by elegant mathematics

physics, chemistry, astronomy, etc.

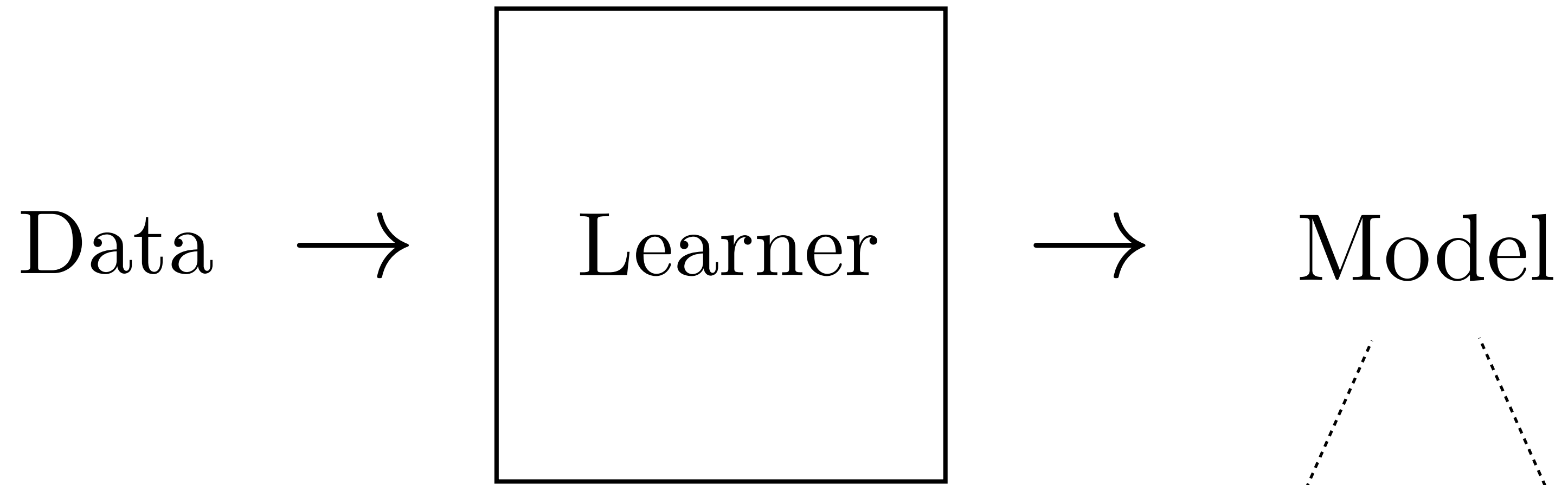
But much cannot

psychology, economics, genetics, etc.

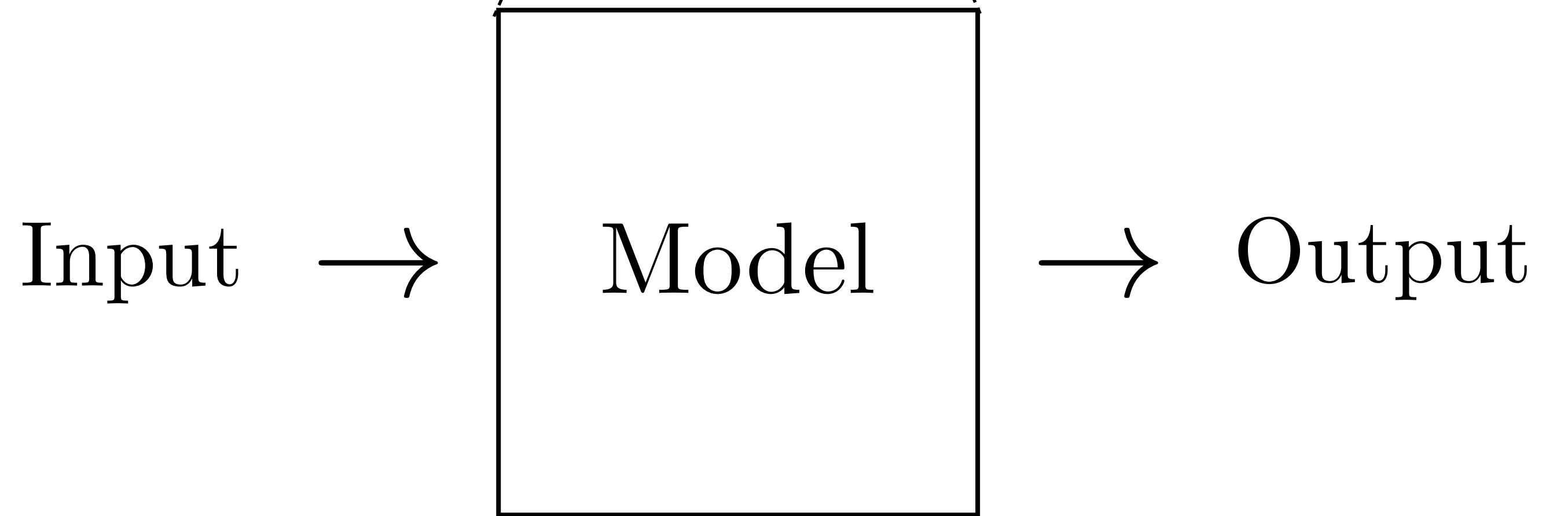
Enter The Data!

“For many tasks, once we have a billion or so examples, we essentially have a closed set that represents (or at least approximates) what we need...”

Learning



Inference



Learning from examples

(aka **supervised learning**)

Training data

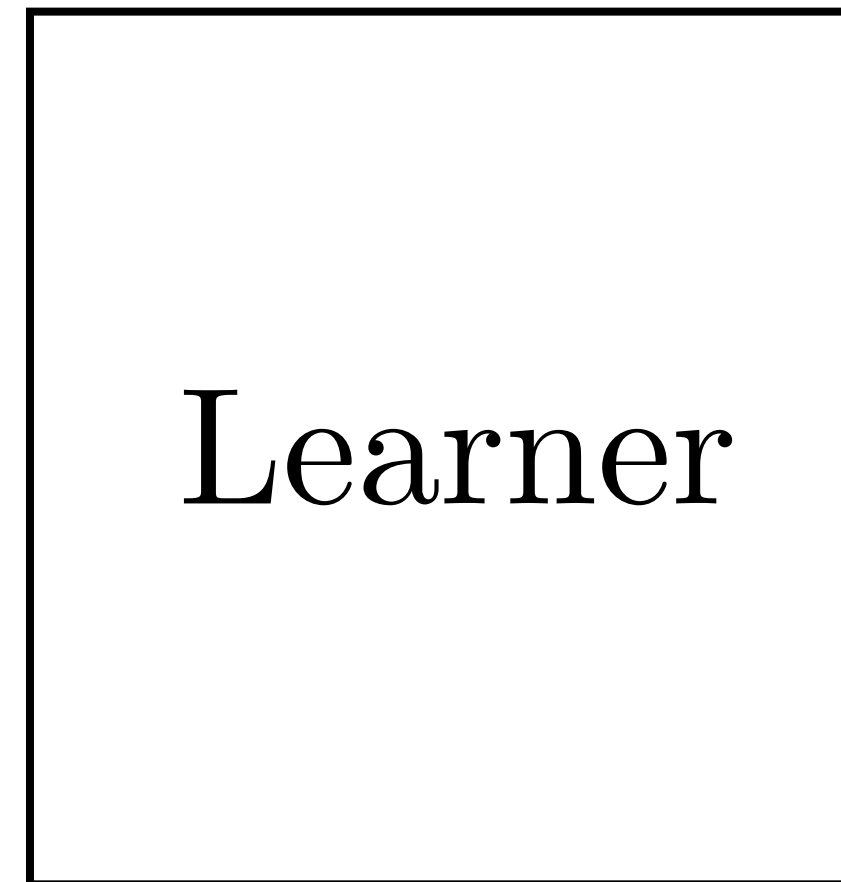
$\{x_1, y_1\}$

$\{x_2, y_2\}$

$\{x_3, y_3\}$

...

→



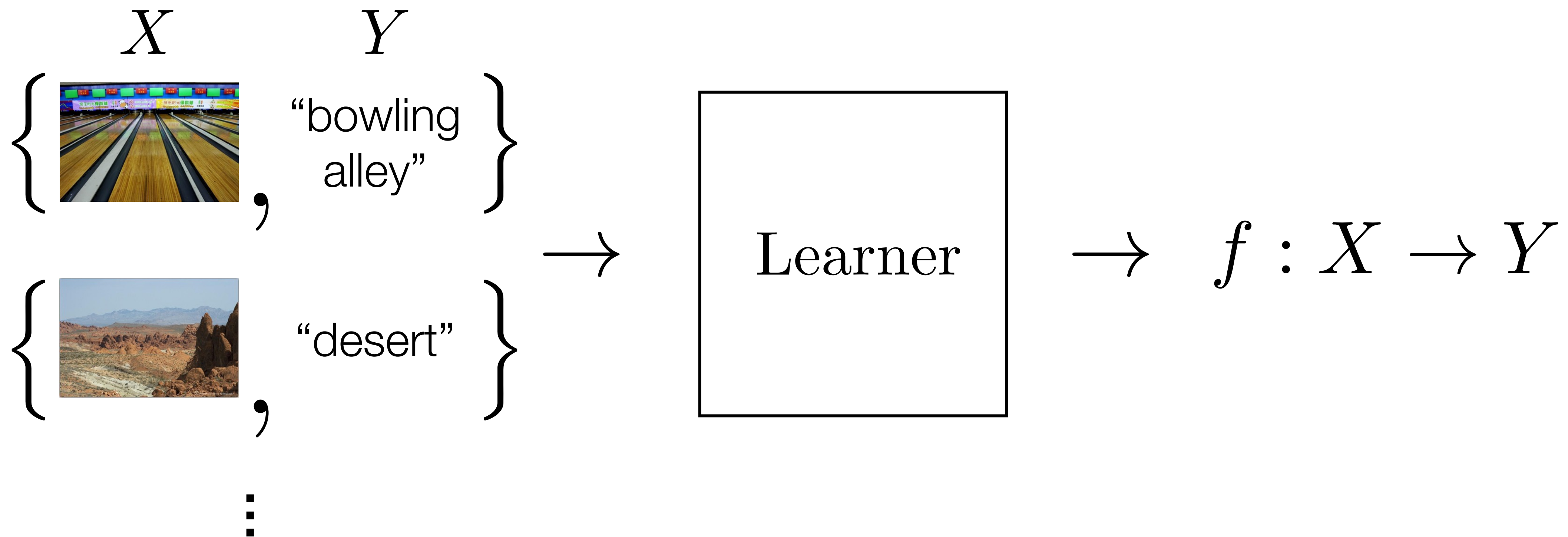
→

$f : X \rightarrow Y$

Learning from examples

(aka **supervised learning**)

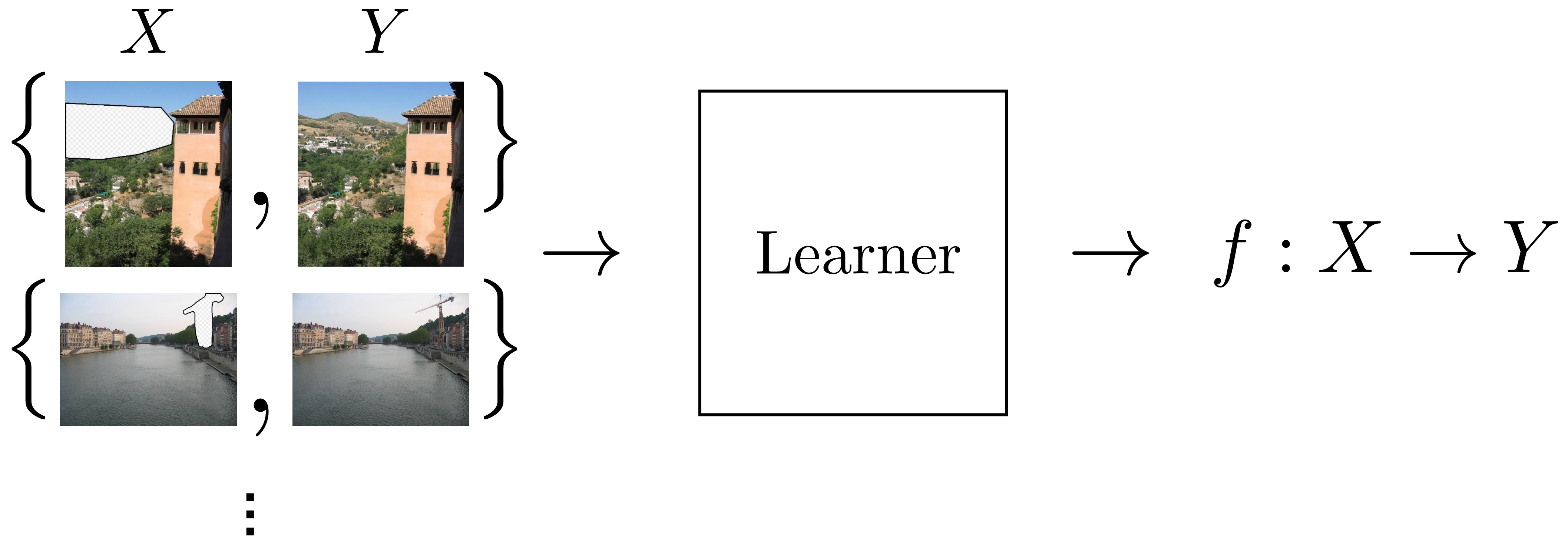
Training data



Learning from examples

(aka **supervised learning**)

Training data



Case study #1: Nearest neighbor

Nearest neighbor

Input feature
vector \mathbf{x}_q



Input image

X



Y

“bowling
alley”

Take label from:
 $\operatorname{argmin}_i \|\mathbf{x}_q - \mathbf{x}_i\|^2$



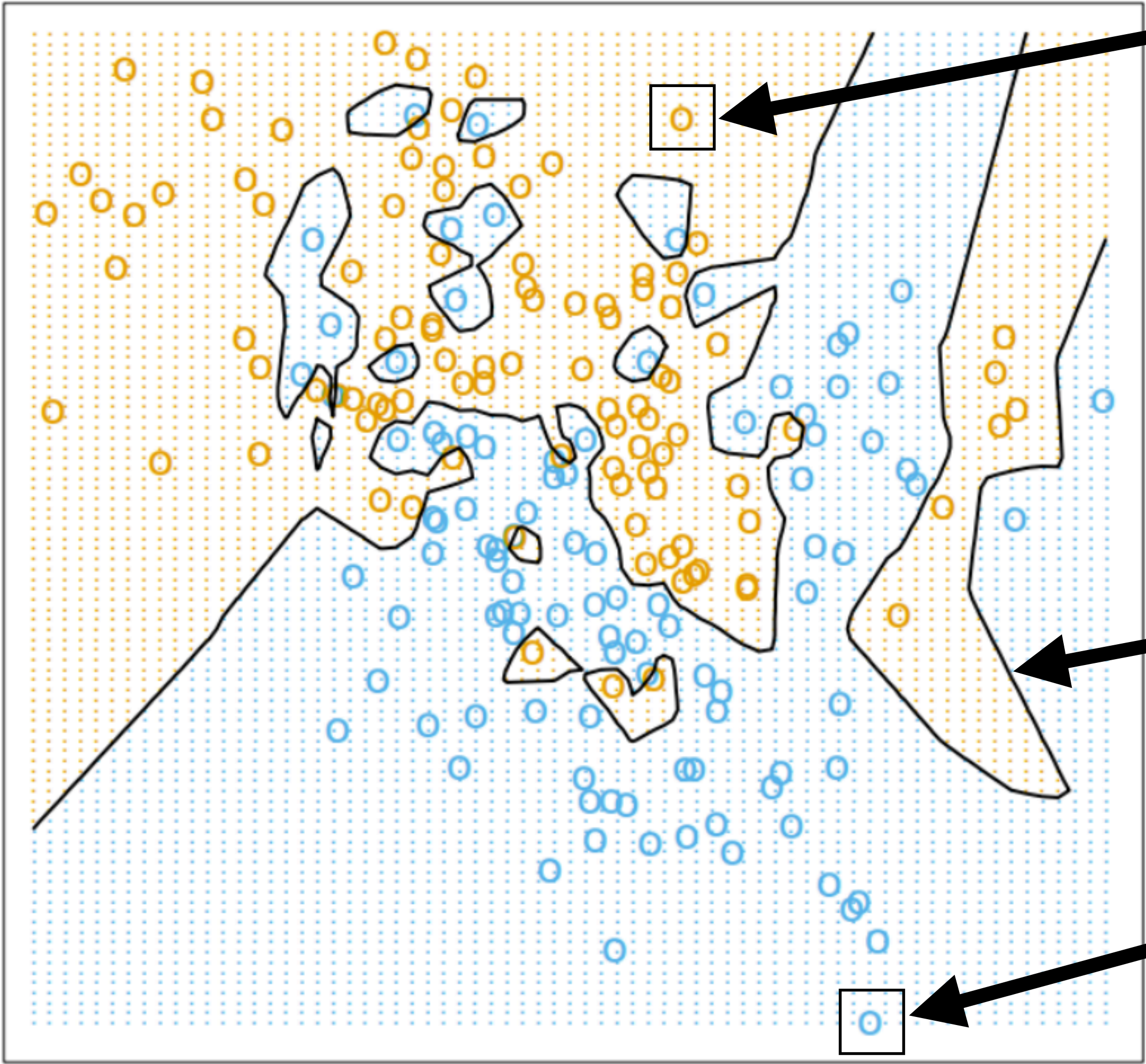
“beach”

K-nearest neighbor:
Take most common
label in k closest
examples



“desert”

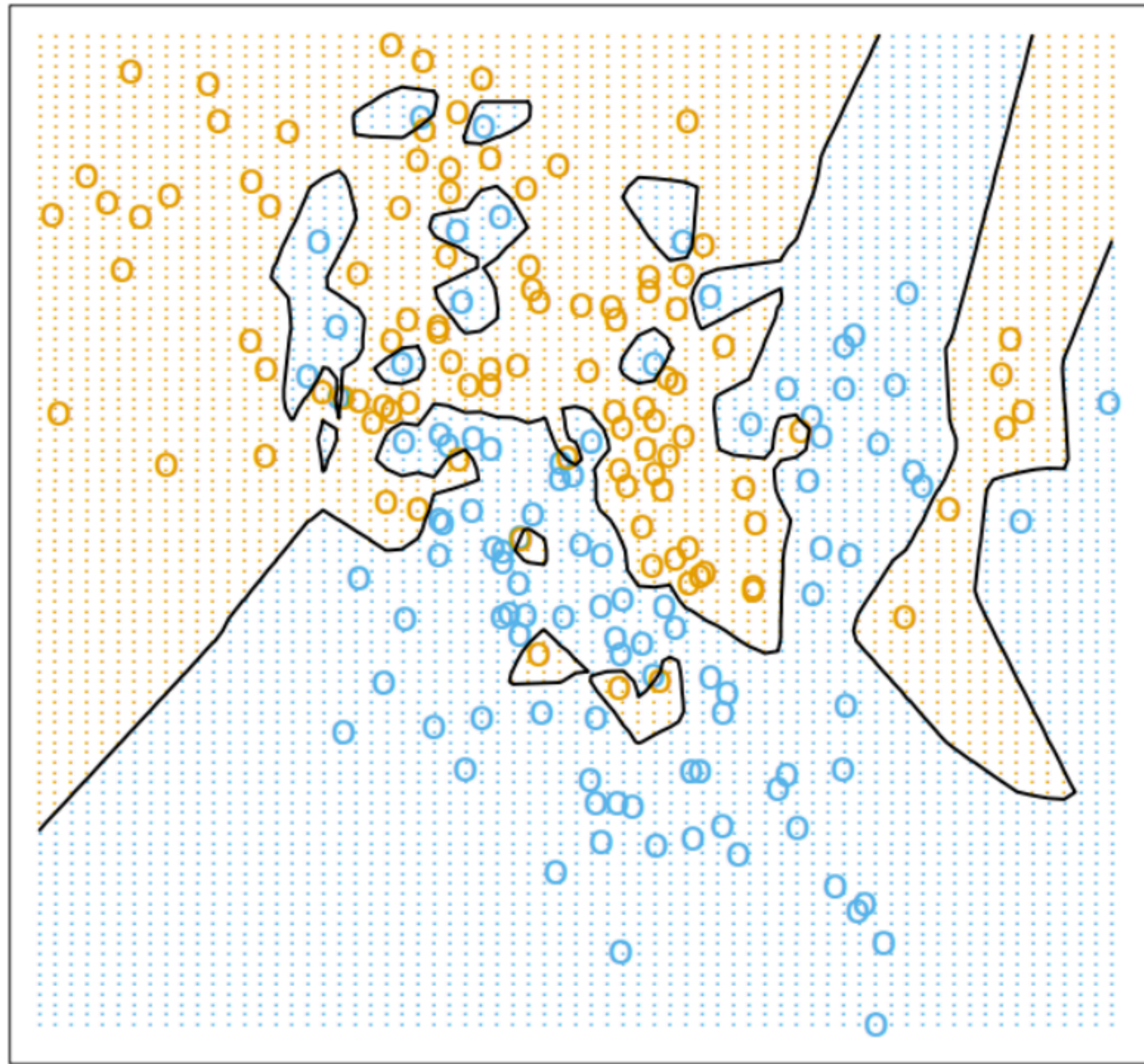
example from orange class



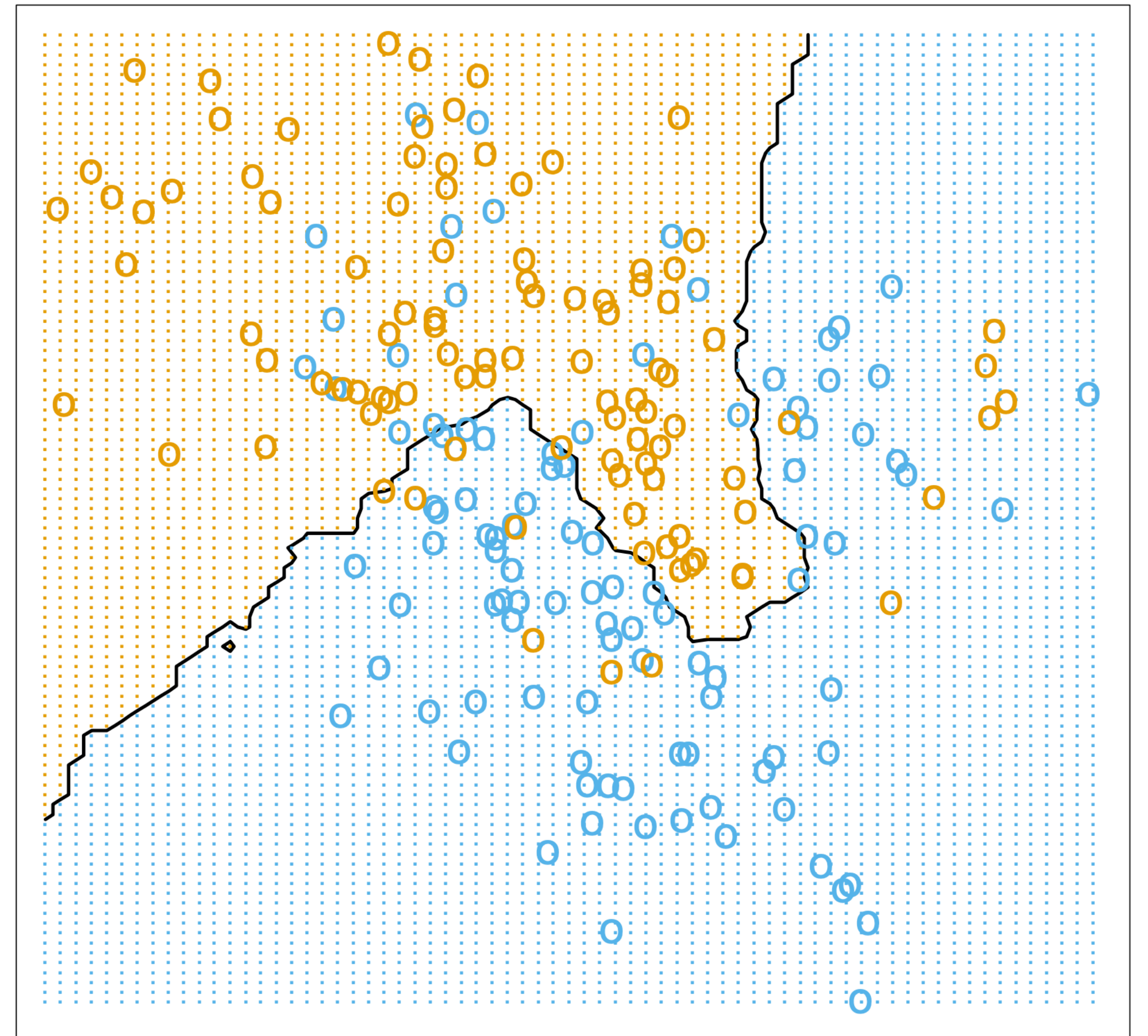
a boundary between
blue/orange classifier predictions

example from blue class

$$k = 1$$



$k = 1$



$k = 15$

How do we choose k ?



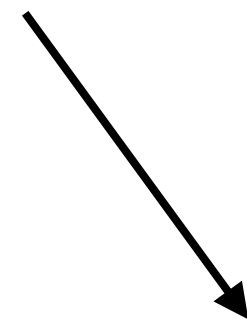
Training Data

The diagram consists of two adjacent rectangular boxes. The left box is light gray and contains the text 'Training Data'. The right box is a darker gray and contains the text 'Test Data'. Both boxes have a thin black border.

Test Data

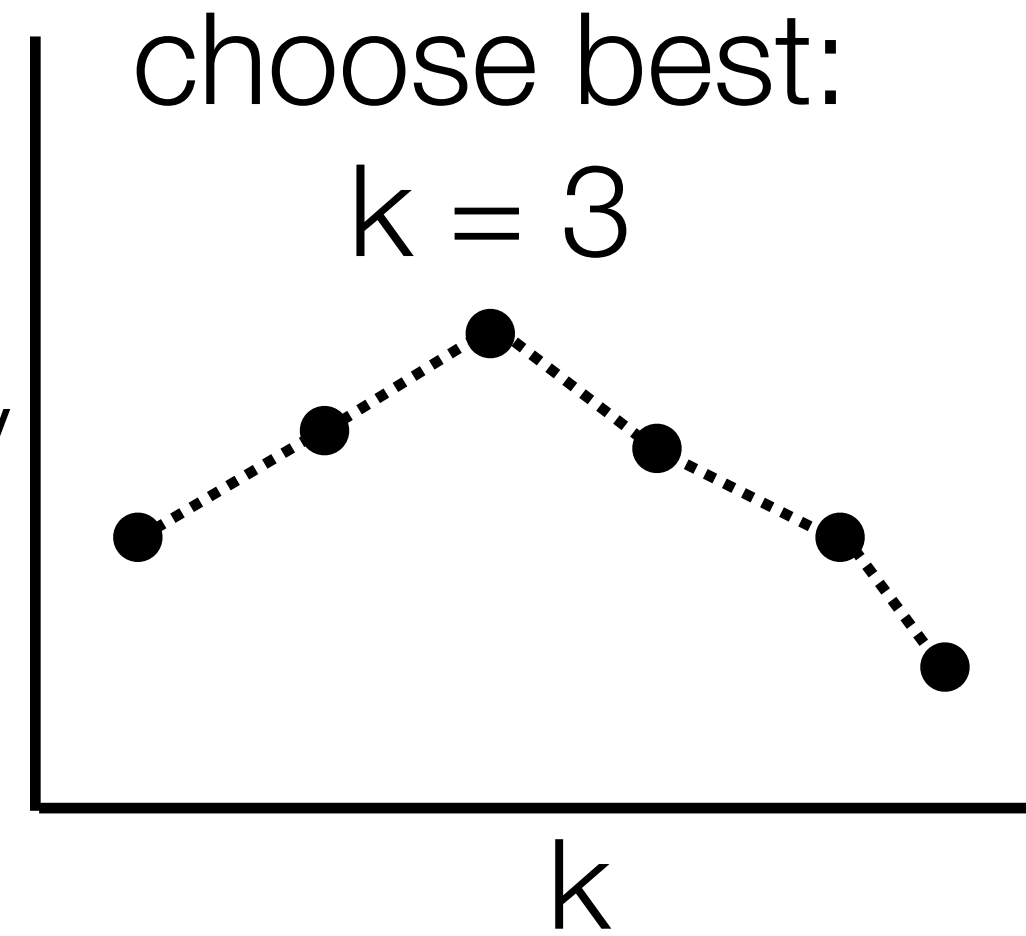
How do we choose k?

Pool of nearest neighbors



Training Set

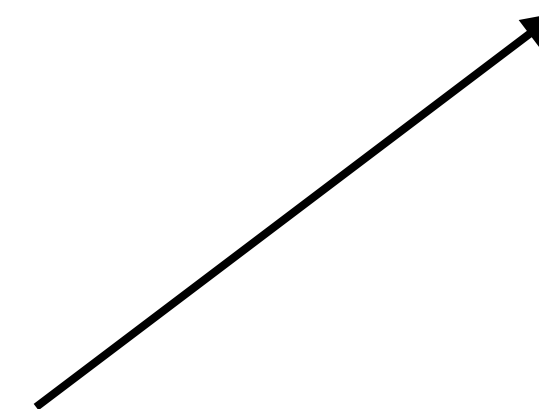
% accuracy



k

Validation Set

Test Set



Measures generalization.
Ideally only test once at end!

Choose **hyperparameters** like k, feature space, similarity function, etc.

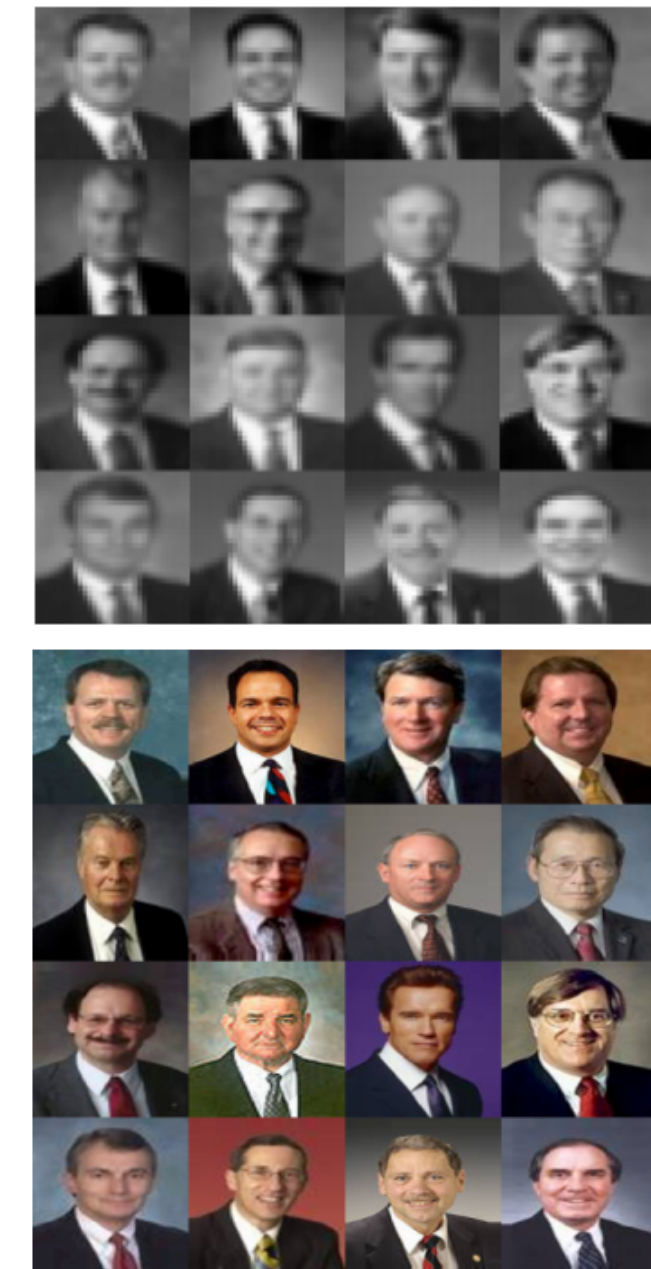
Can work well when datasets are very big!



Input image



Colorized output



Nearest Neighbors

[Torralba et al. "80 million tiny images", 2008]

How do you represent the image?

Idea: L2 distance between pixels



Spot the difference



Not Blurred



Blurred

Spot the difference



Spot the difference



Shifted

Spot the difference



Original

L2 distance: 51.1



Shifted

L2 distance: 51.6



Blurred

Almost the same distance!

How do you represent the image?

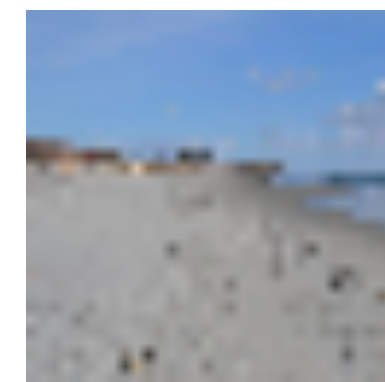


Input image

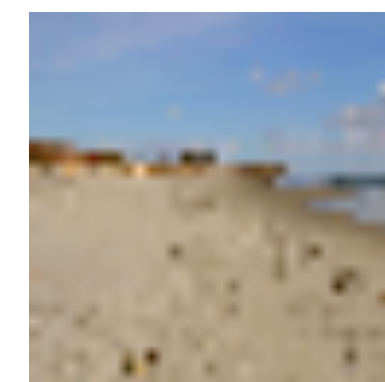
→
$$\begin{bmatrix} 0.5 \\ 0.1 \\ -0.4 \\ \dots \\ 0.9 \end{bmatrix}$$

Feature vector x_q

Some feature ideas:

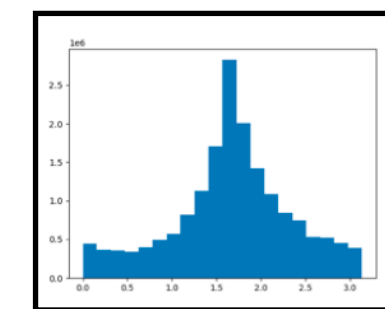


tiny image



normalized tiny image

$$\frac{x - \mu}{\sigma}$$

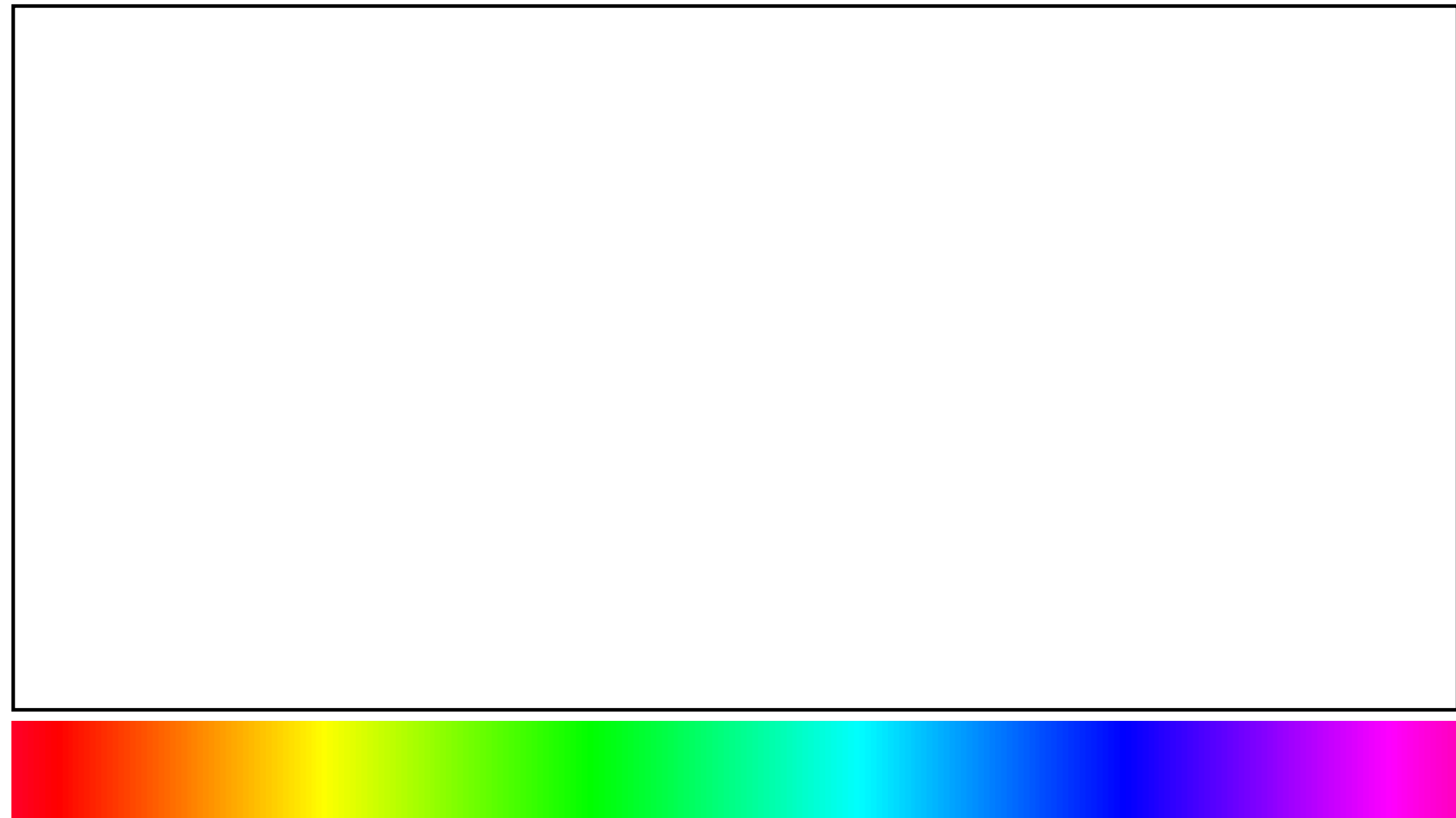


edge orientation histogram

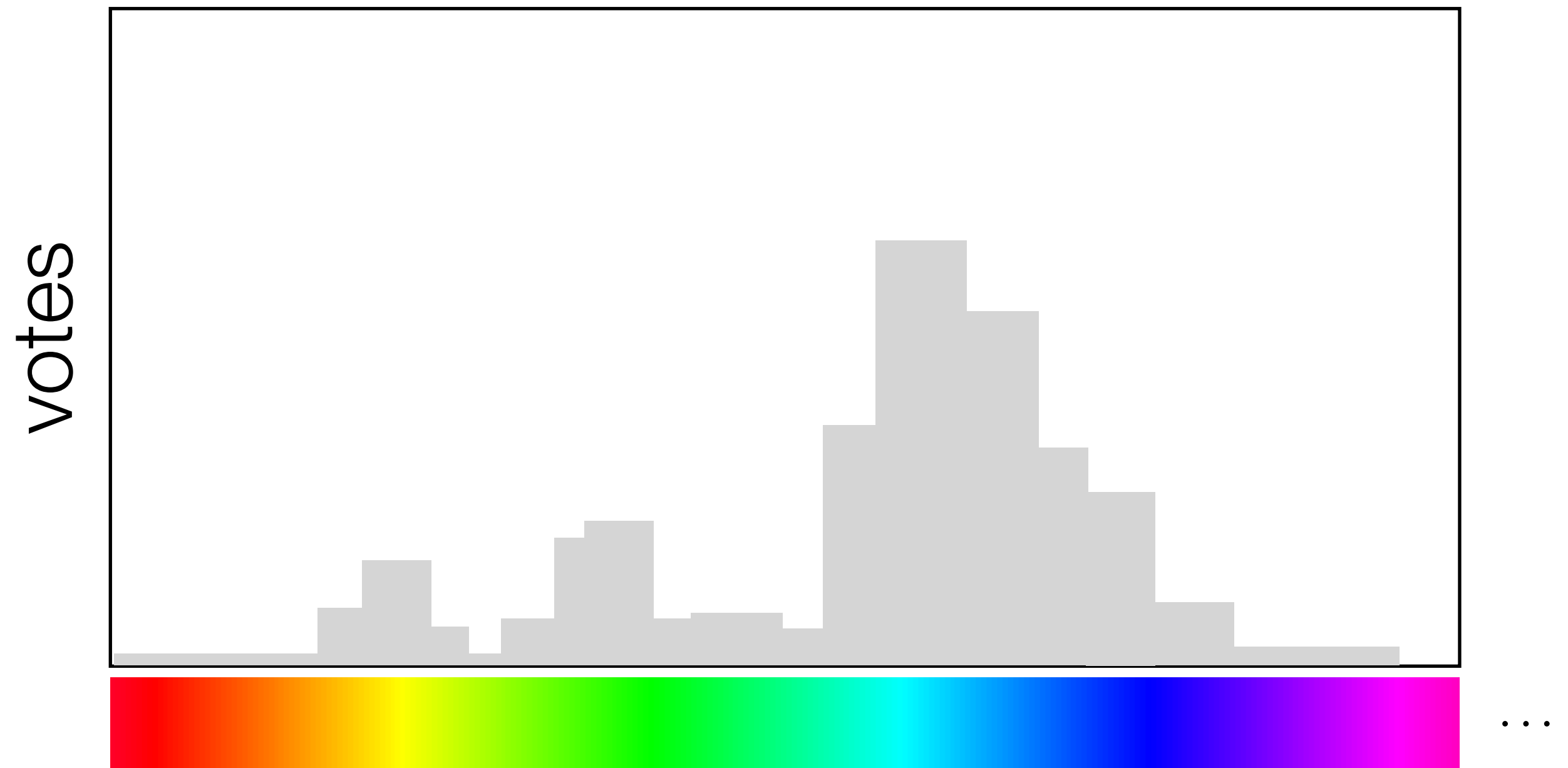
Color histograms



votes

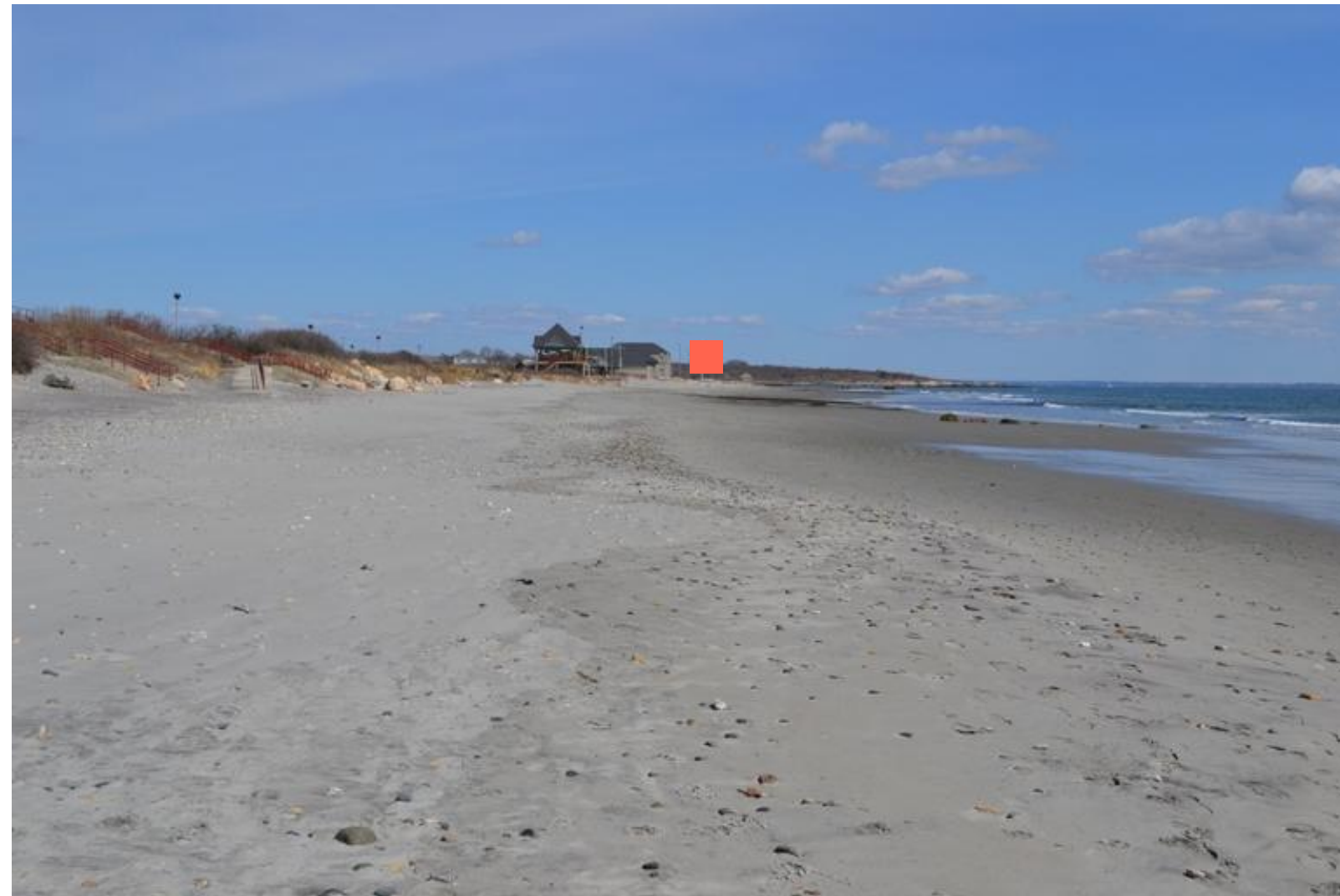


Color histograms



Typically create a coarse grid over the color channels, such as $10 \times 10 \times 10$.

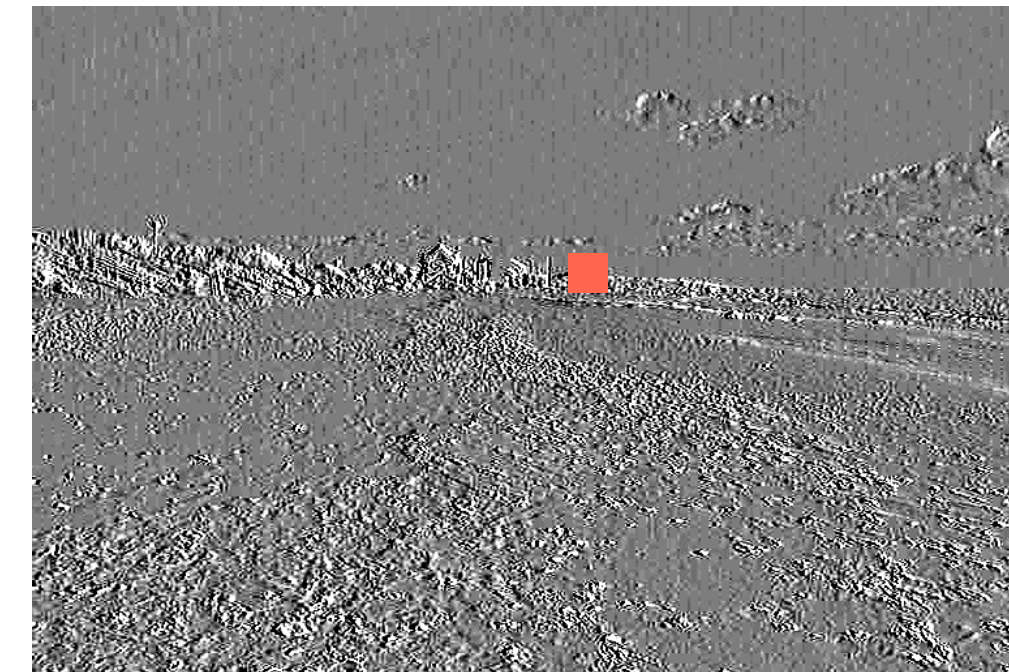
Edge orientation histograms



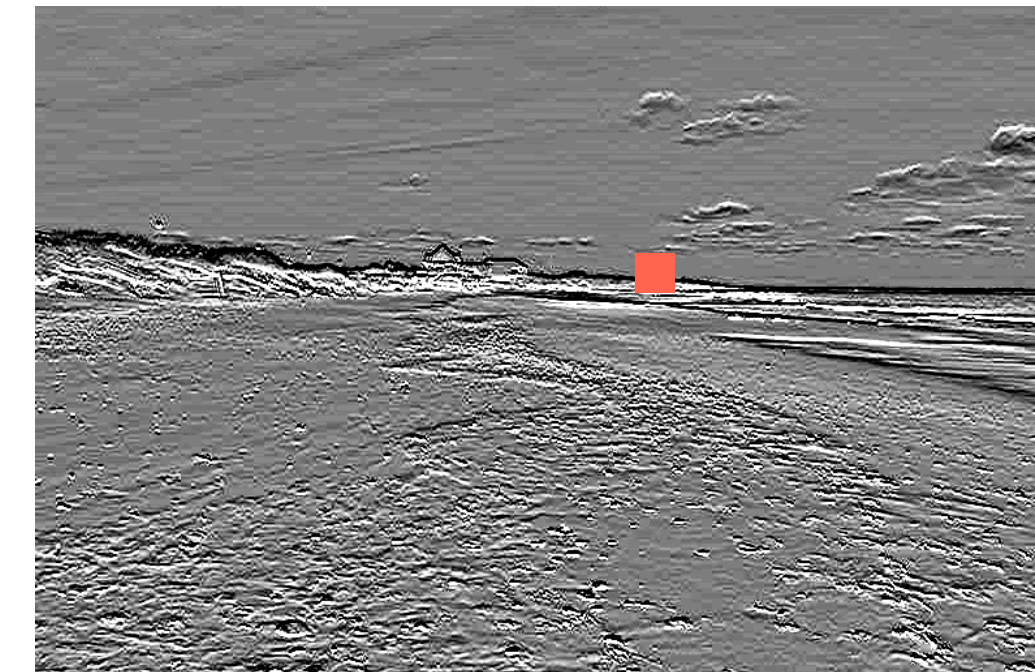
Input image



magnitude

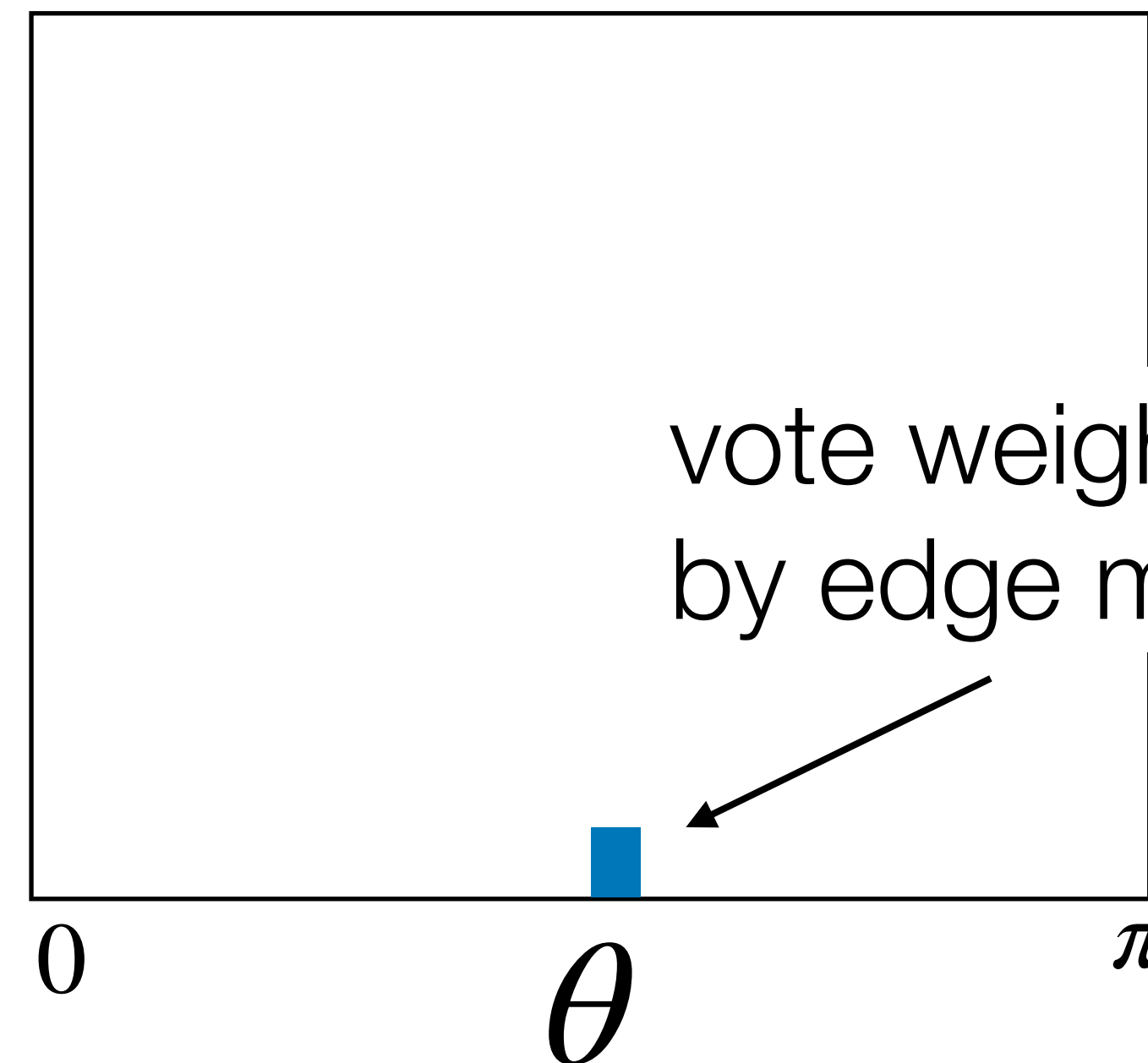


dx



dy

$$\theta = \tan^{-1}(dy/dx)$$



vote weighted
by edge magnitude

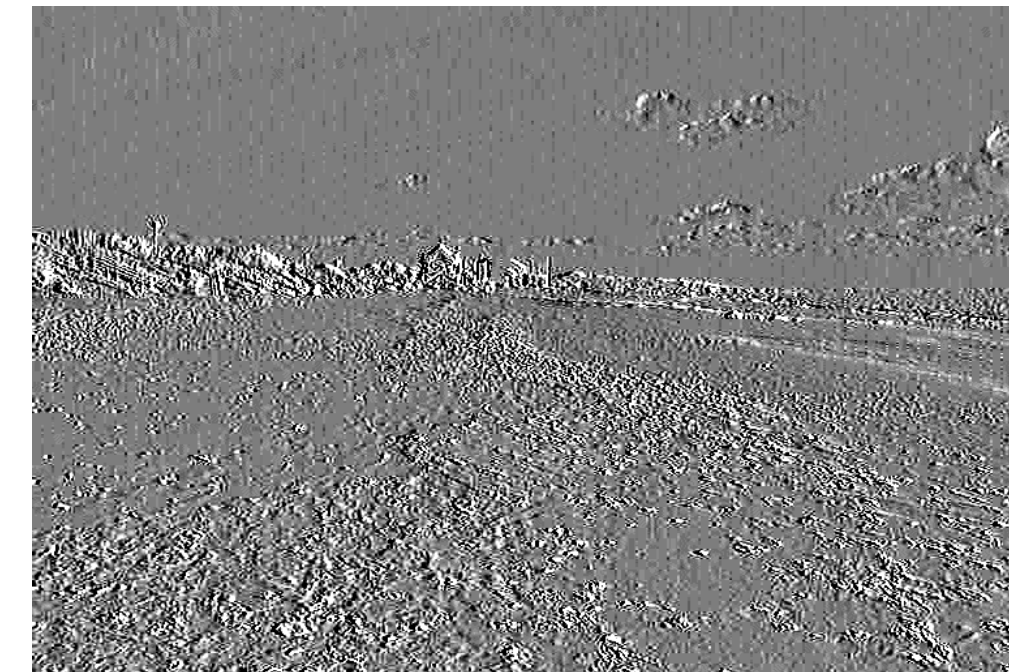
Edge orientation histograms



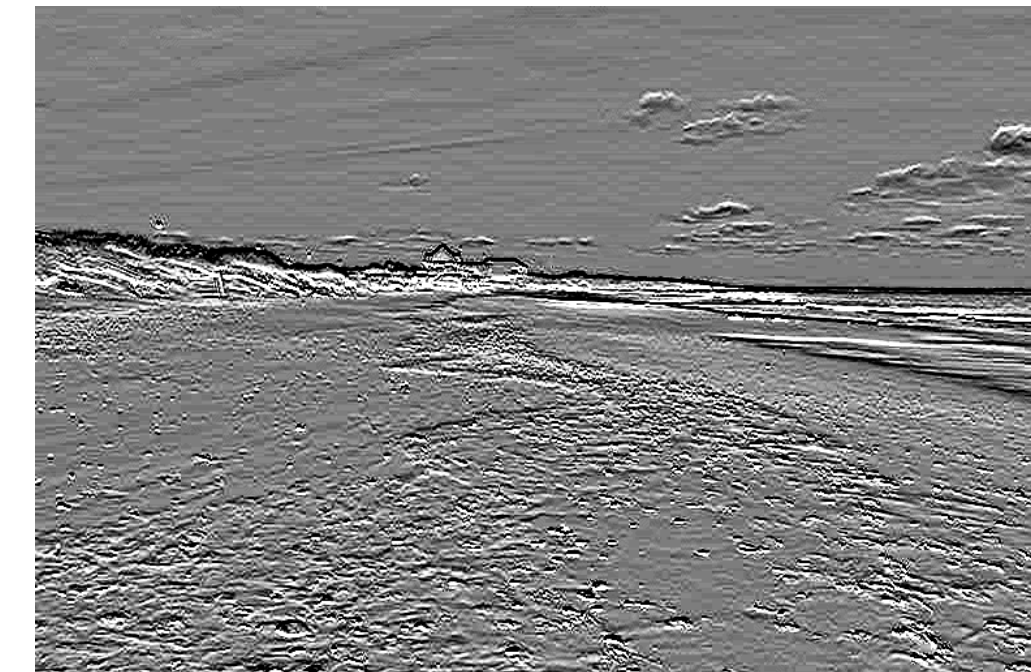
Input image



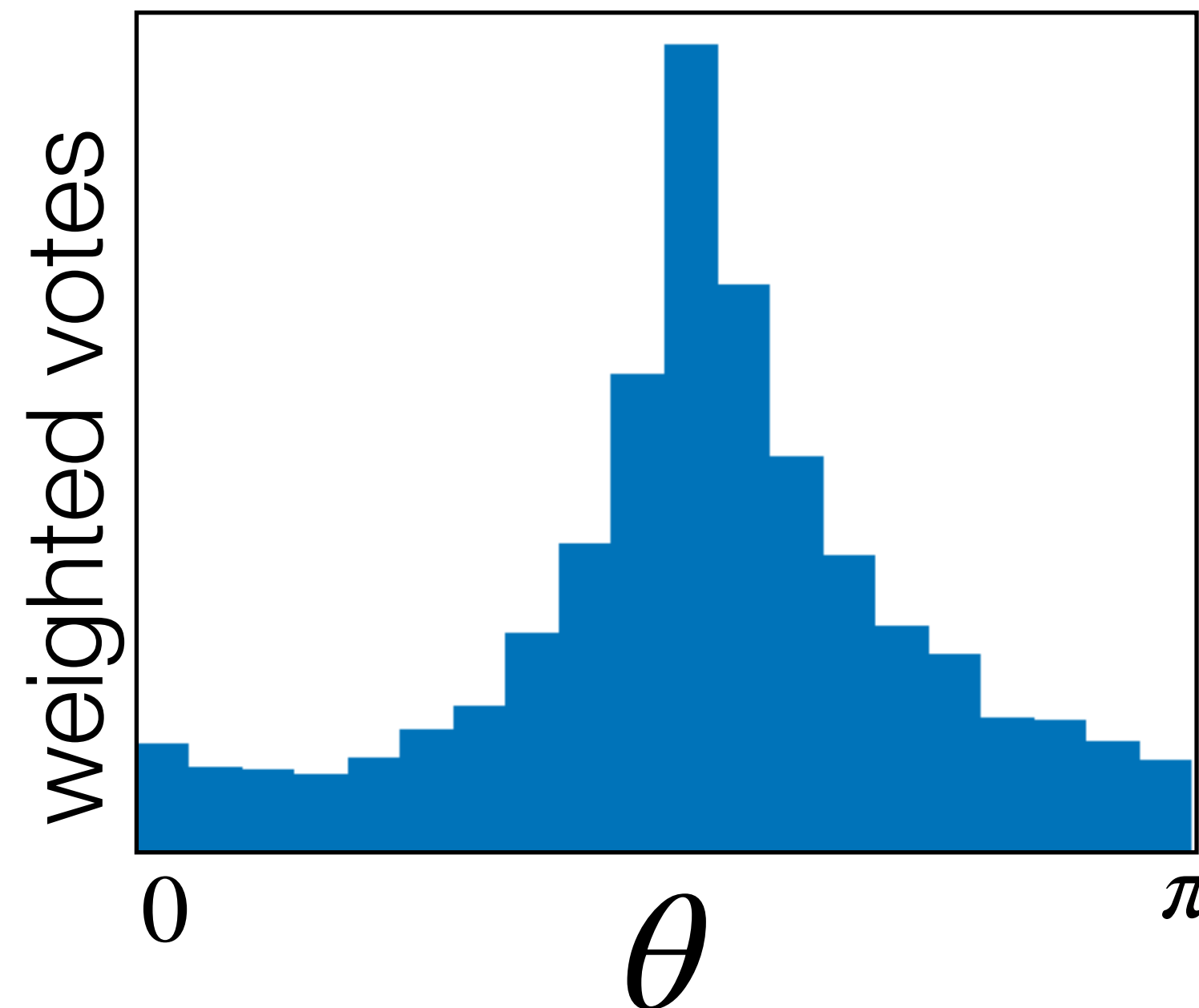
magnitude



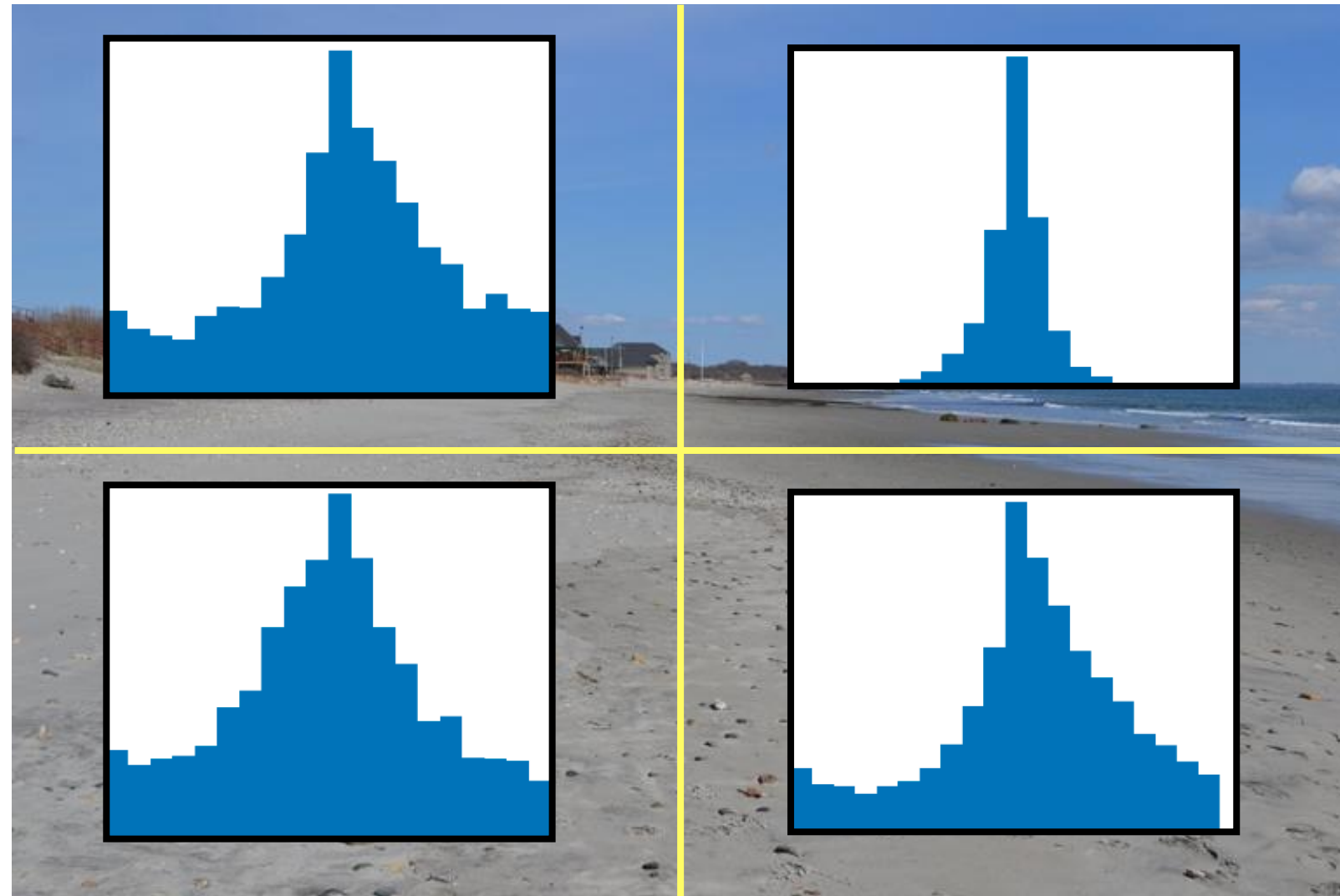
dx



dy



Edge orientation histograms

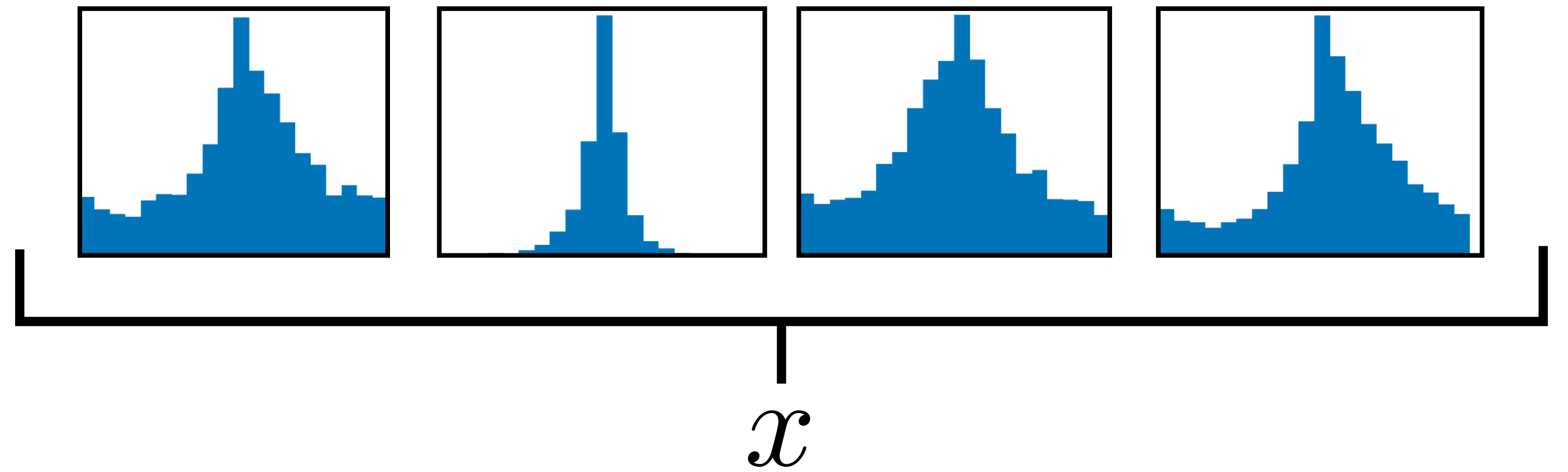


Input image

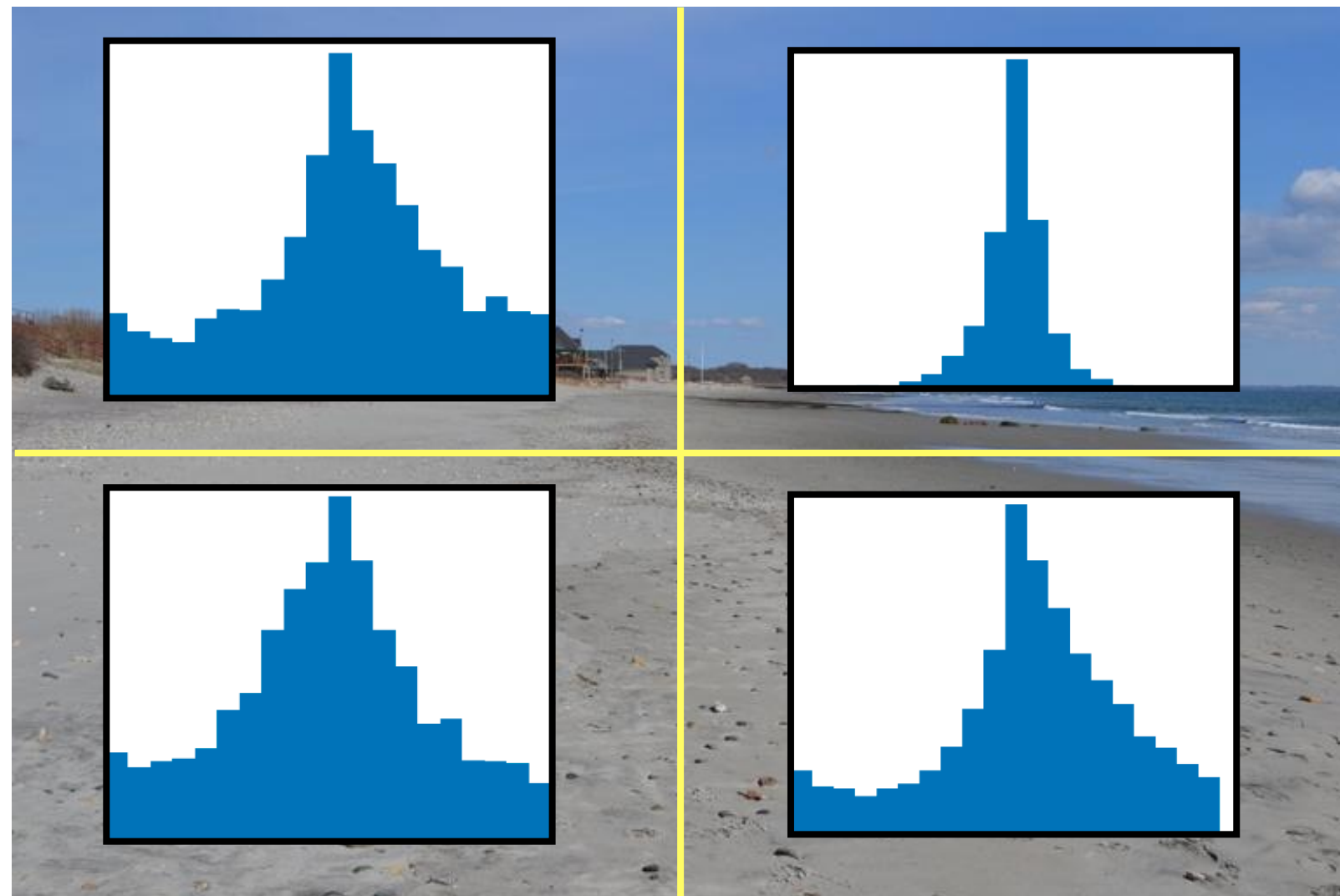
Edge orientation histograms



Input image



Edge orientation histograms

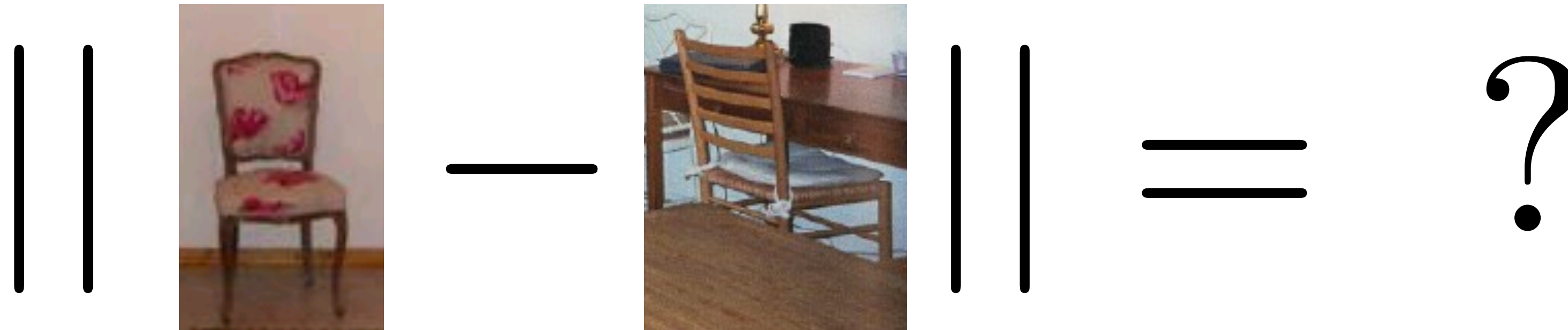


Input image

- Common variation is called: Histograms of Oriented Gradients (HOG) [Dalal and Triggs 2005].
- Often these use more complex normalization and weighting schemes.
- Used for image matching with SIFT features [Lowe 1999].
- We'll discuss *feature learning* methods next week. These usually perform a lot better!

Nearest Neighbor

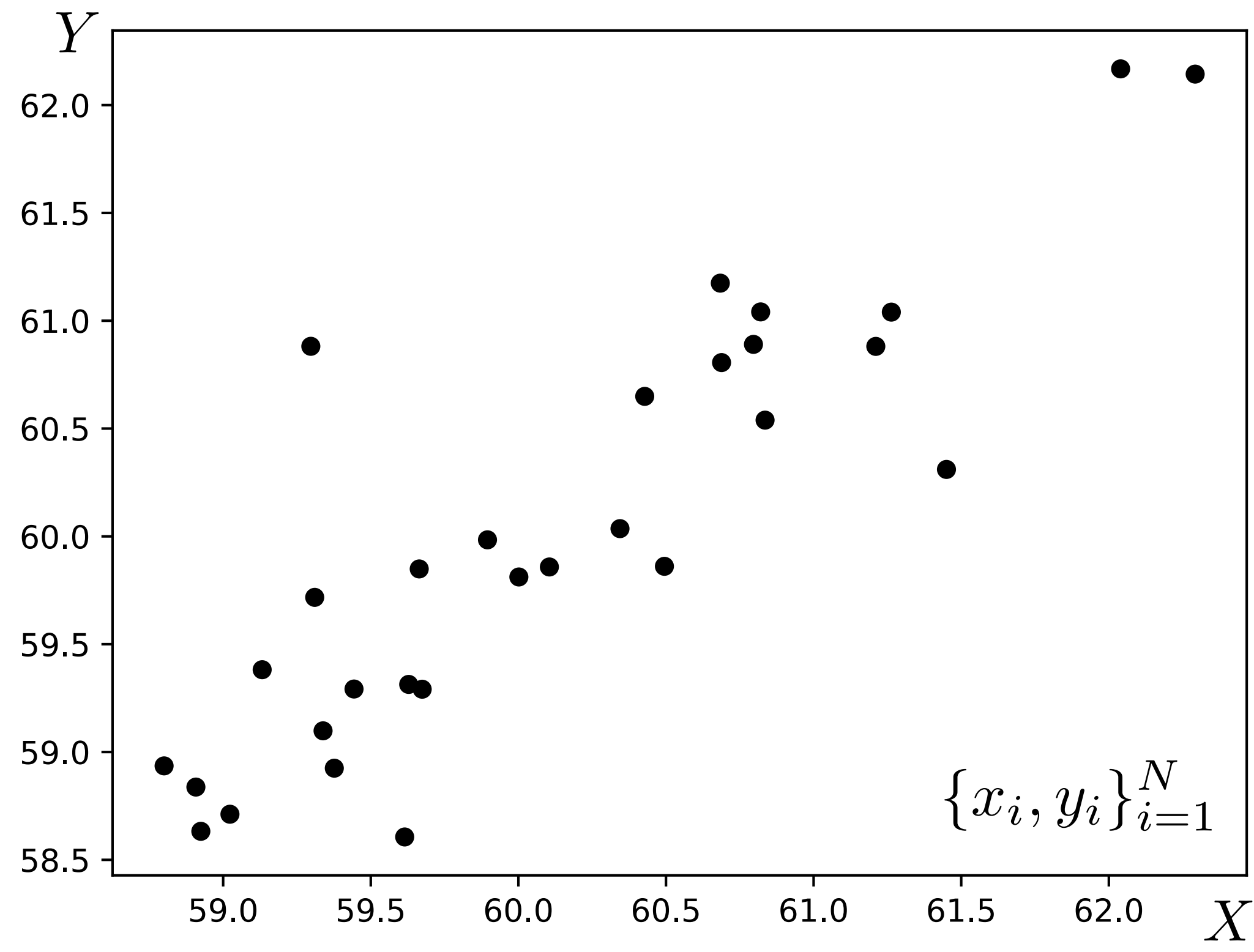
- Hard to define similarity metric



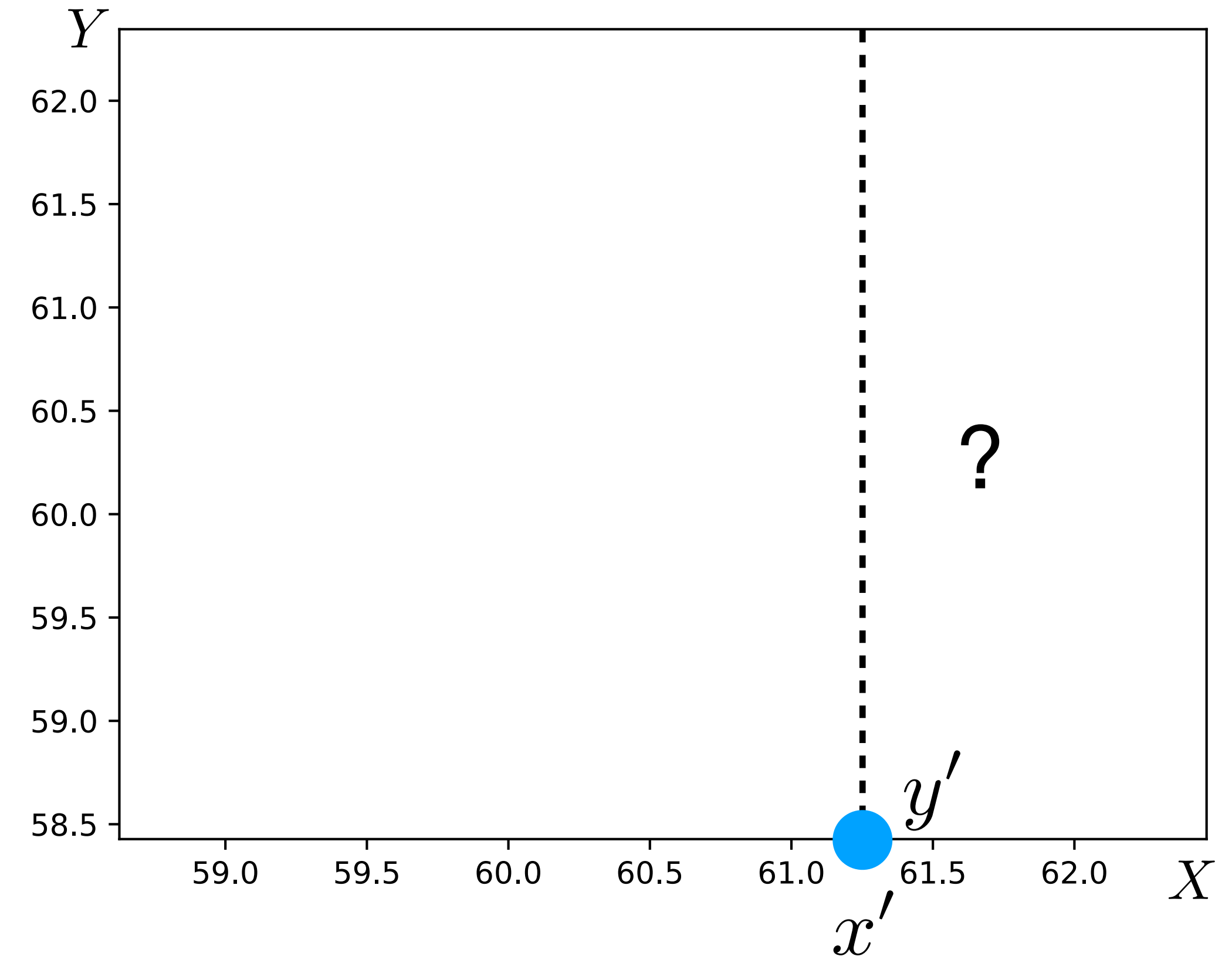
- Pays attention to all features equally
- Can be very slow: $O(nd)$

Case study #2: Linear least squares

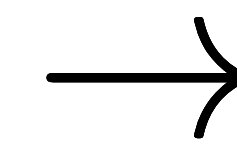
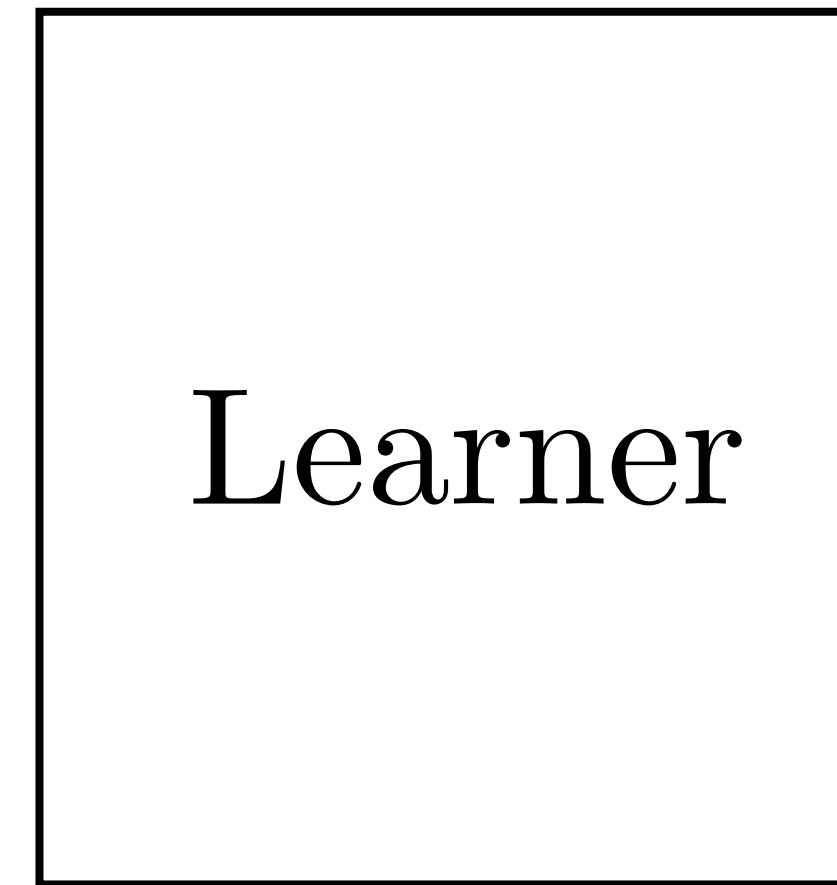
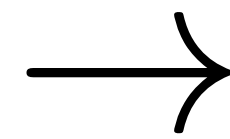
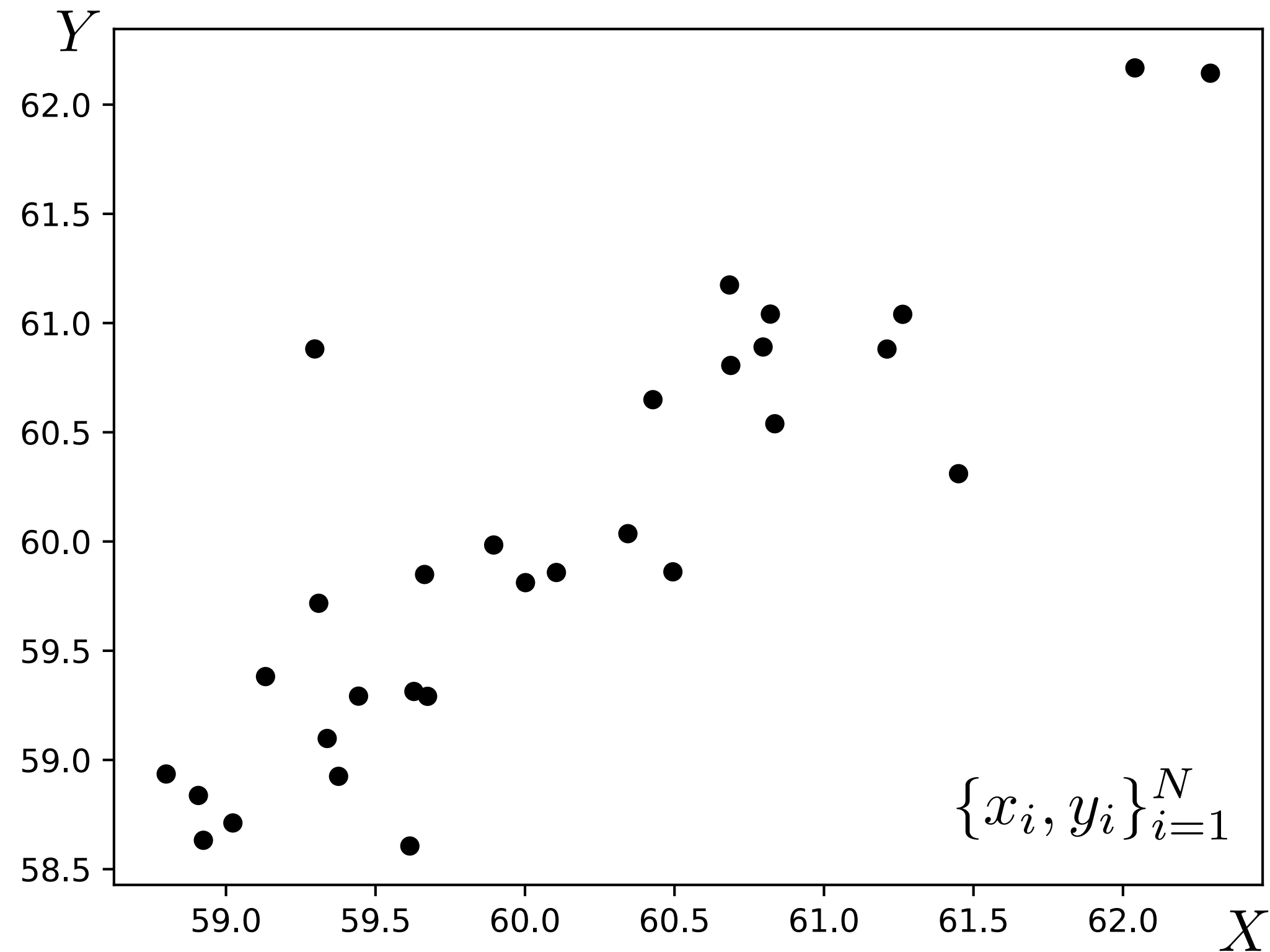
Training data



Test query



Training data

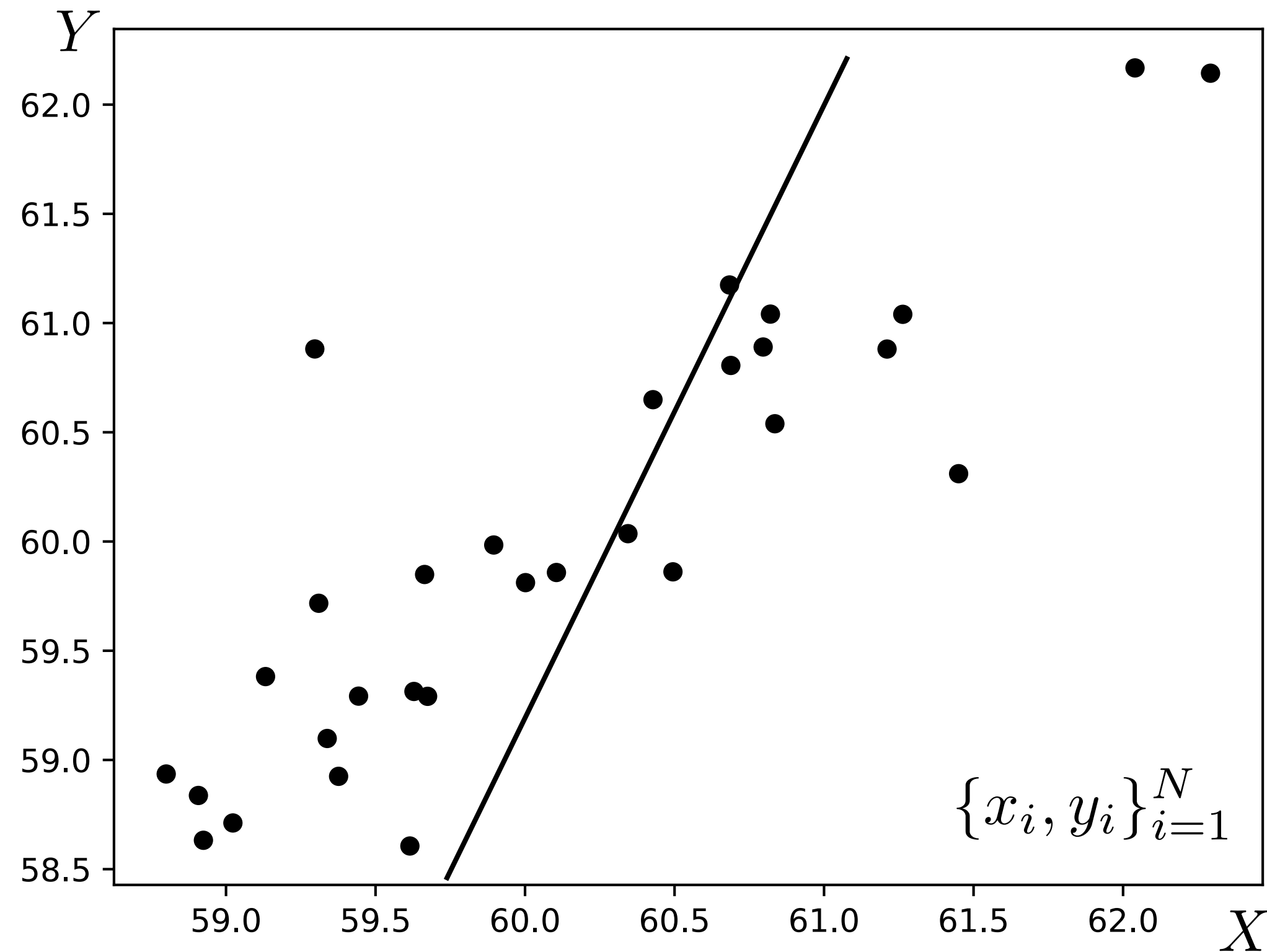


$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Hypothesis space

The relationship between X and Y is roughly linear: $y \approx \theta_1 x + \theta_0$

Training data

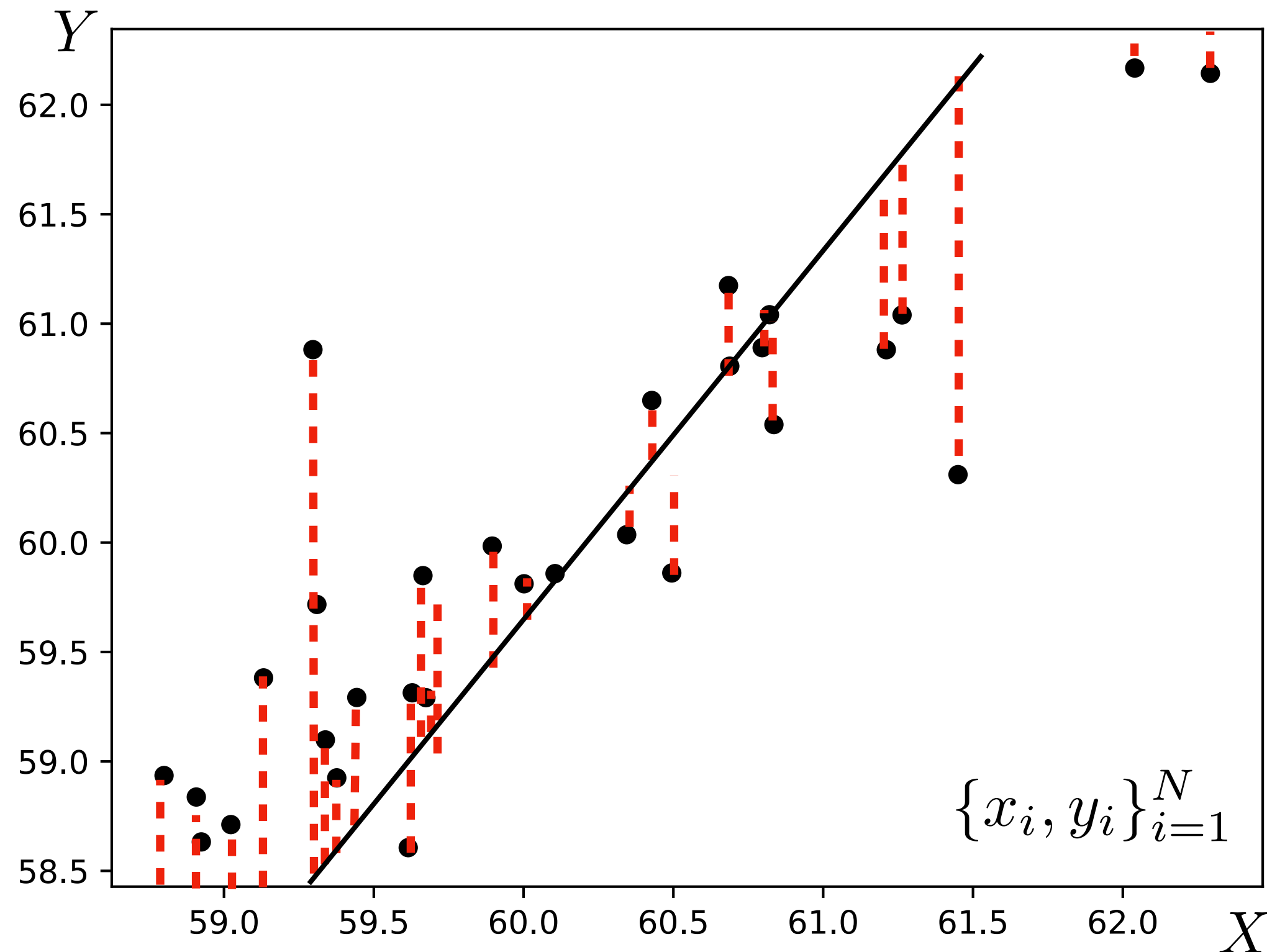


Search for the **parameters**, $\theta = \{\theta_0, \theta_1\}$, that best fit the data.

$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Best fit in what sense?

Training data



Search for the **parameters**, $\theta = \{\theta_0, \theta_1\}$, that best fit the data.

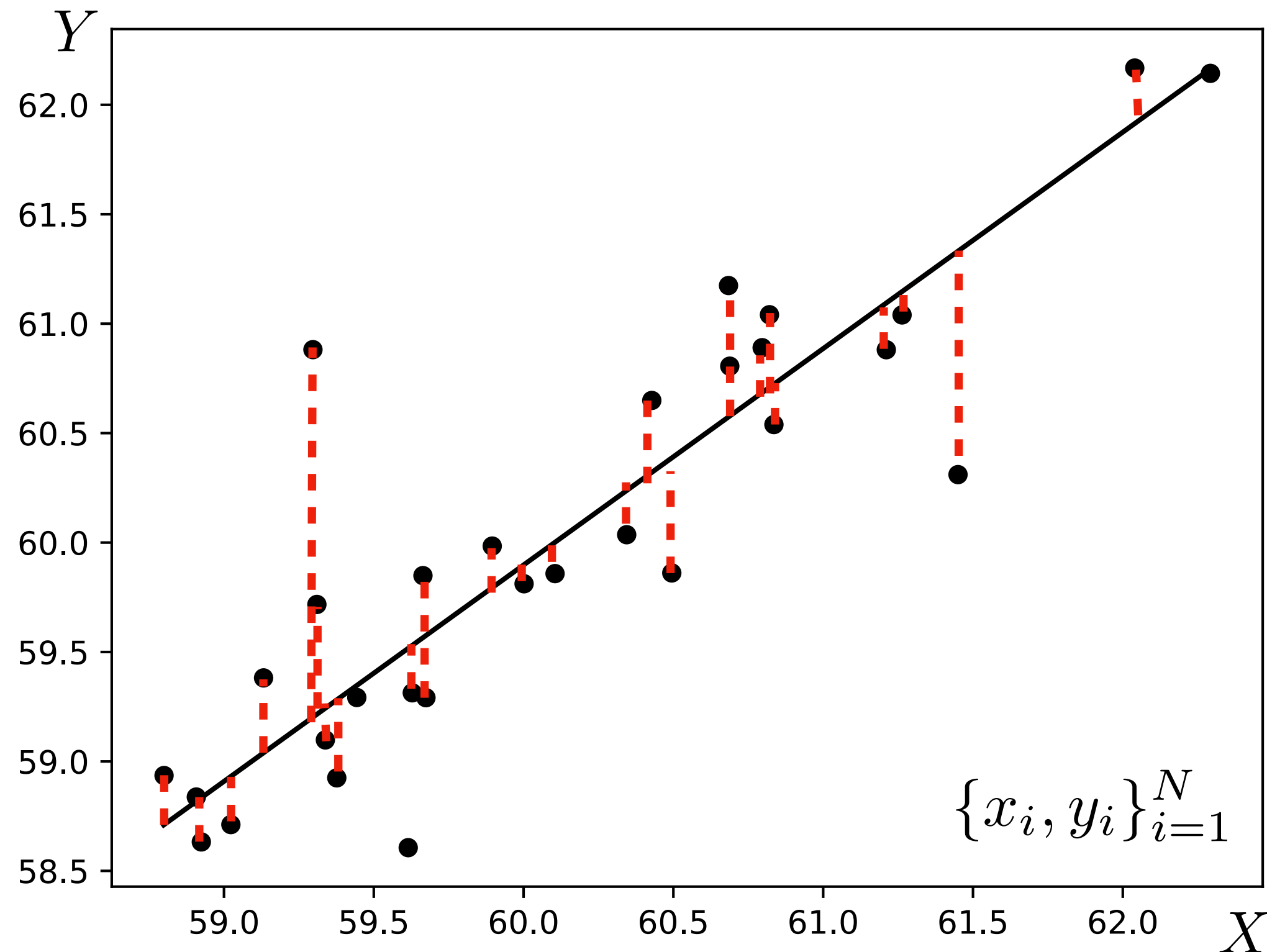
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Best fit in what sense?

The least-squares **objective** (aka **loss**) says the best fit is the function that minimizes the squared error between predictions and target values:

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2 \quad \hat{y} \equiv f_{\theta}(x)$$

Training data



Search for the **parameters**, $\theta = \{\theta_0, \theta_1\}$, that best fit the data.

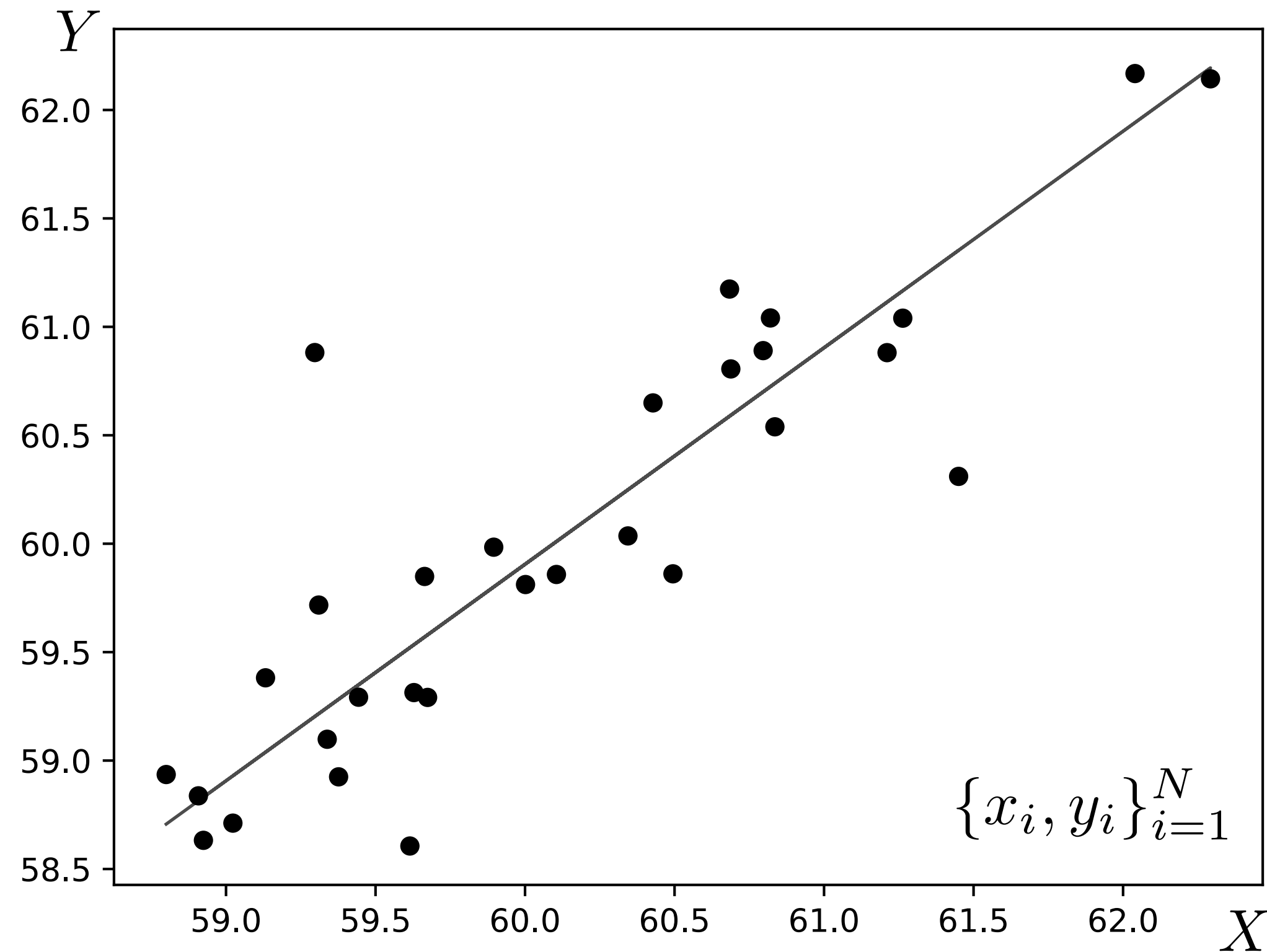
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Best fit in what sense?

The least-squares **objective** (aka **loss**) says the best fit is the function that minimizes the squared error between predictions and target values:

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2 \quad \hat{y} \equiv f_{\theta}(x)$$

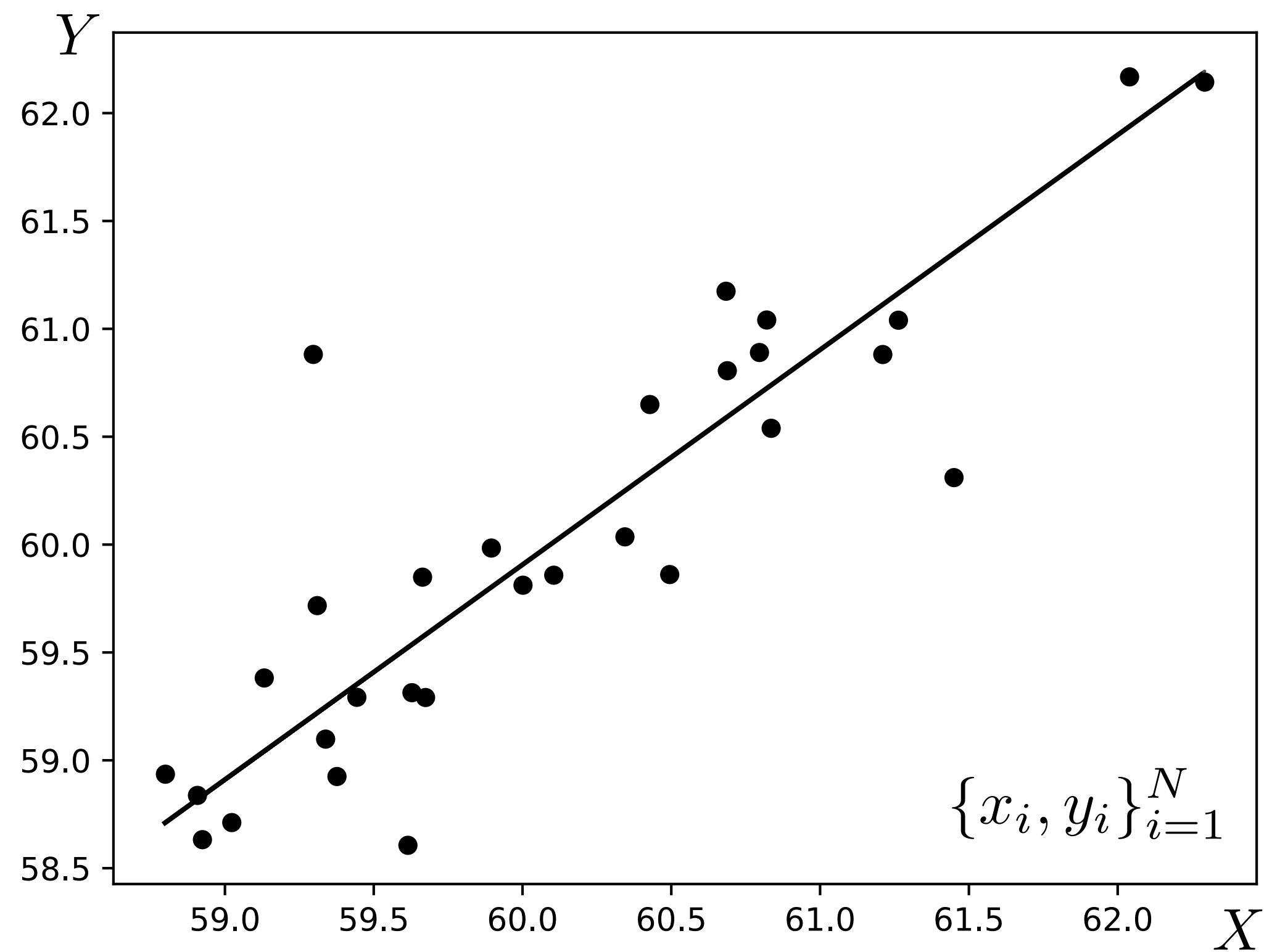
Training data



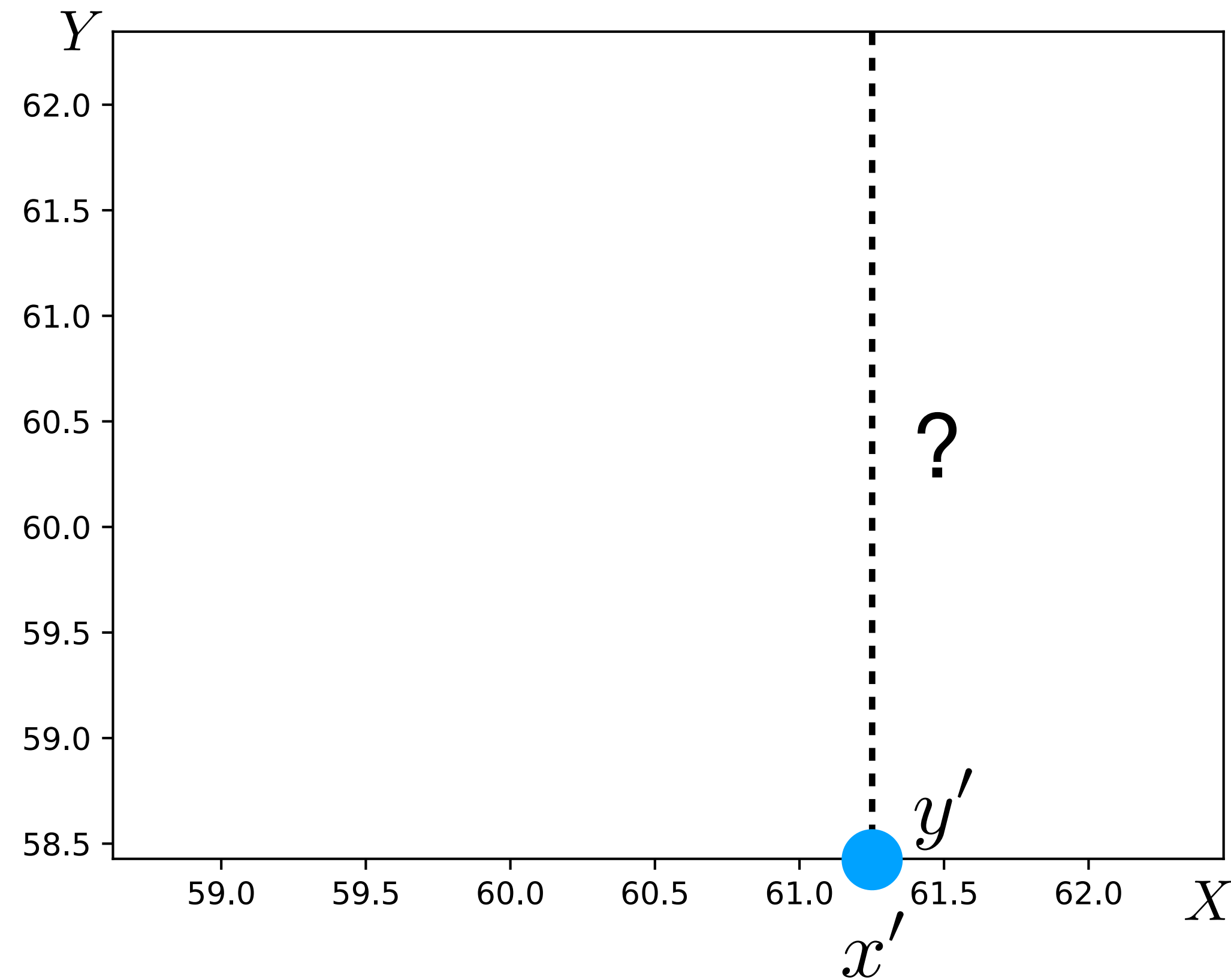
Complete learning problem:

$$\begin{aligned}\theta^* &= \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x_i) - y_i)^2 \\ &= \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x_i + \theta_0 - y_i)^2\end{aligned}$$

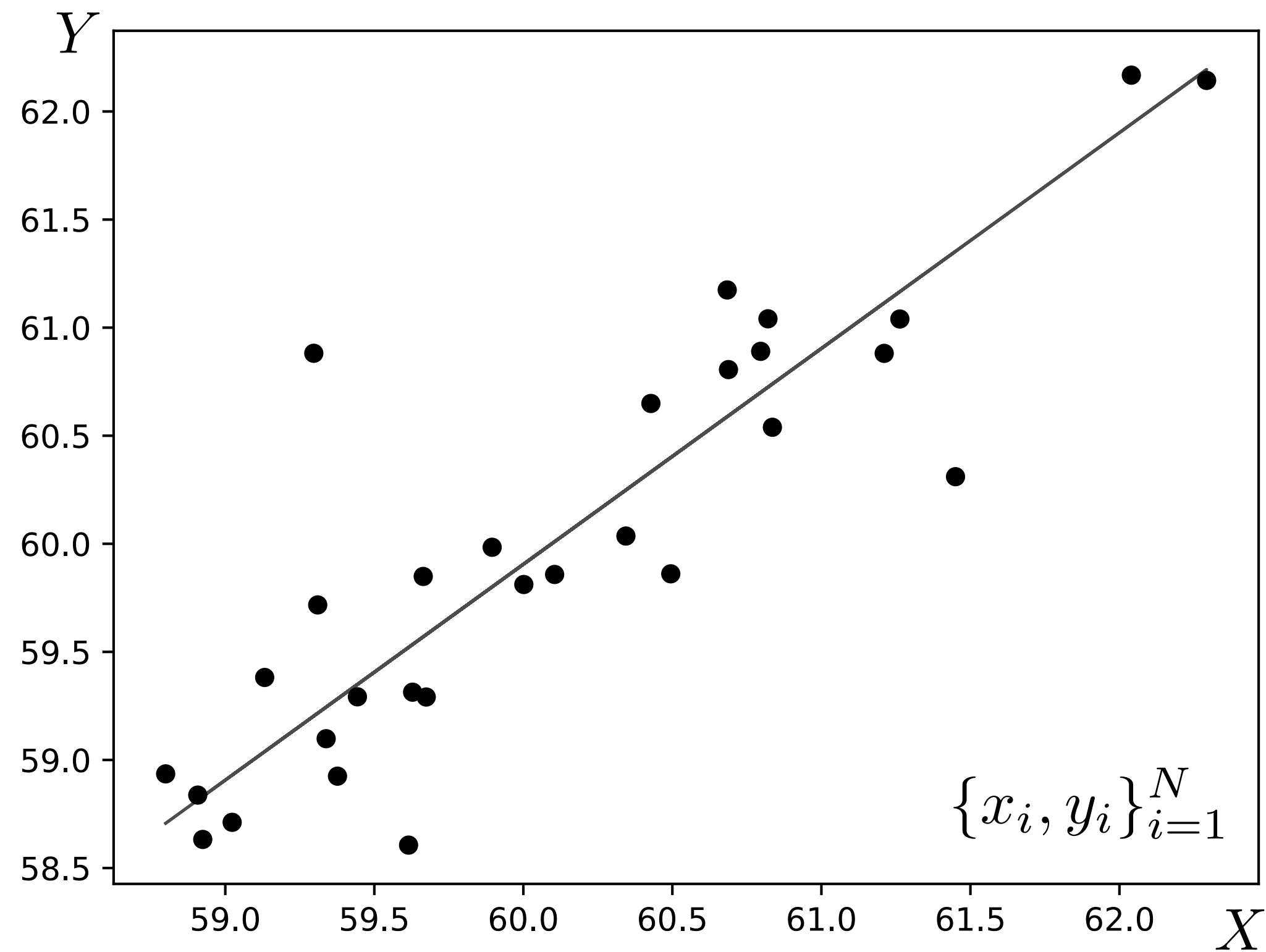
Training data



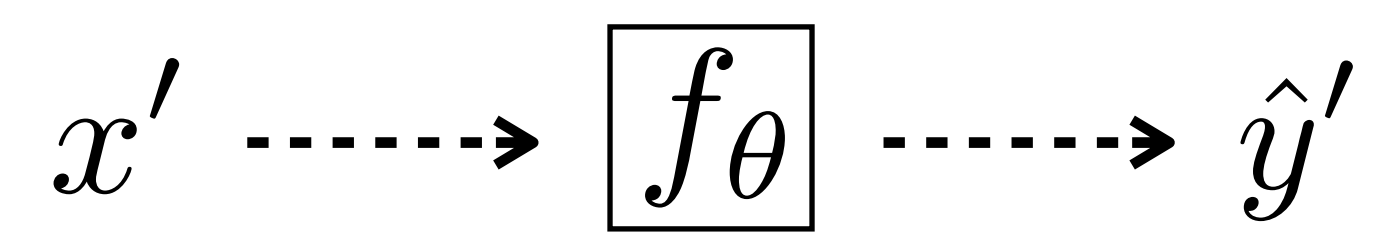
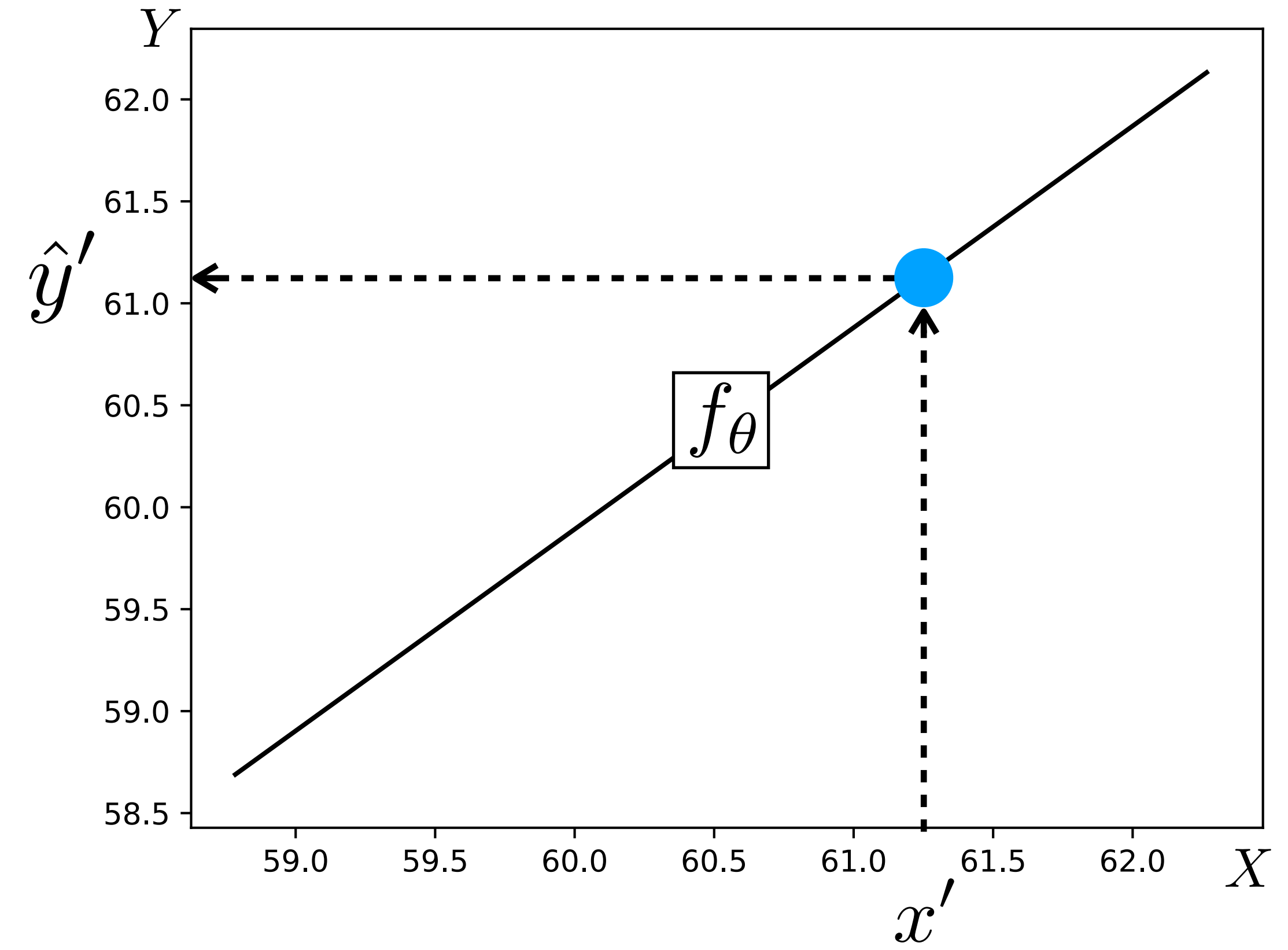
Test query



Training data



Test query



How to minimize the objective w.r.t. θ ?

Learning problem

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x_i + \theta_0 - y_i)^2$$

Recall:

$$\begin{aligned} J(\theta) &= \sum_{I=1}^N (\theta_1 x_i + \theta_0 - y_i)^2 \\ &= (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \end{aligned}$$

$$\mathbf{X} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_N & 1 \end{pmatrix} \quad \theta = (\theta_1 \quad \theta_0) \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

$$\theta^* = \arg \min_{\theta} J(\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta} = 0$$

$$\frac{\partial J(\theta)}{\partial \theta} = 2(\mathbf{X}^T \mathbf{X}\theta - \mathbf{X}^T \mathbf{y})$$

$$2(\mathbf{X}^T \mathbf{X}\theta^* - \mathbf{X}^T \mathbf{y}) = 0$$

$$\mathbf{X}^T \mathbf{X}\theta^* = \mathbf{X}^T \mathbf{y}$$

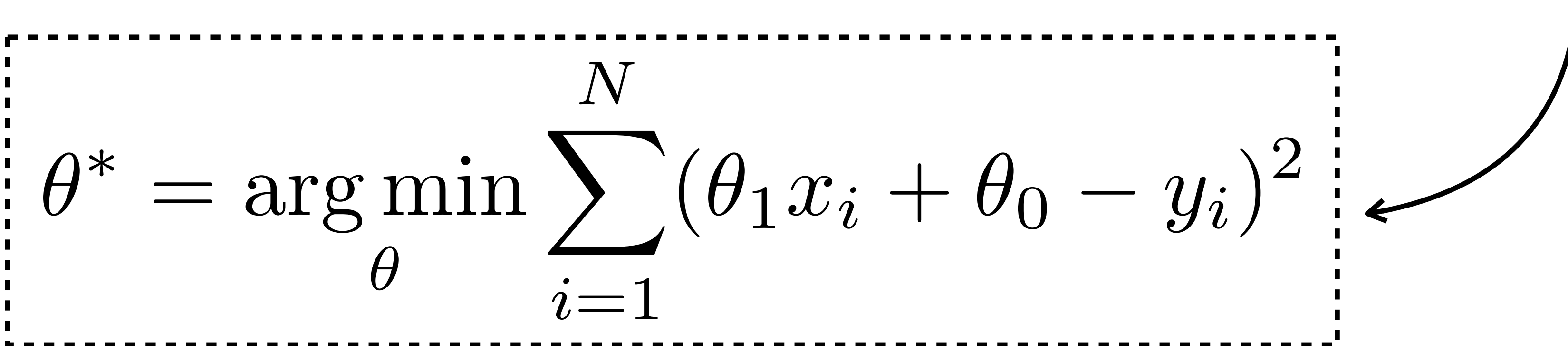
$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Solution

Empirical Risk Minimization

(formalization of supervised learning)

Linear least squares learning problem

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x_i + \theta_0 - y_i)^2$$


Empirical Risk Minimization

(formalization of supervised learning)

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i), \mathbf{y}_i)$$

Objective function (loss)

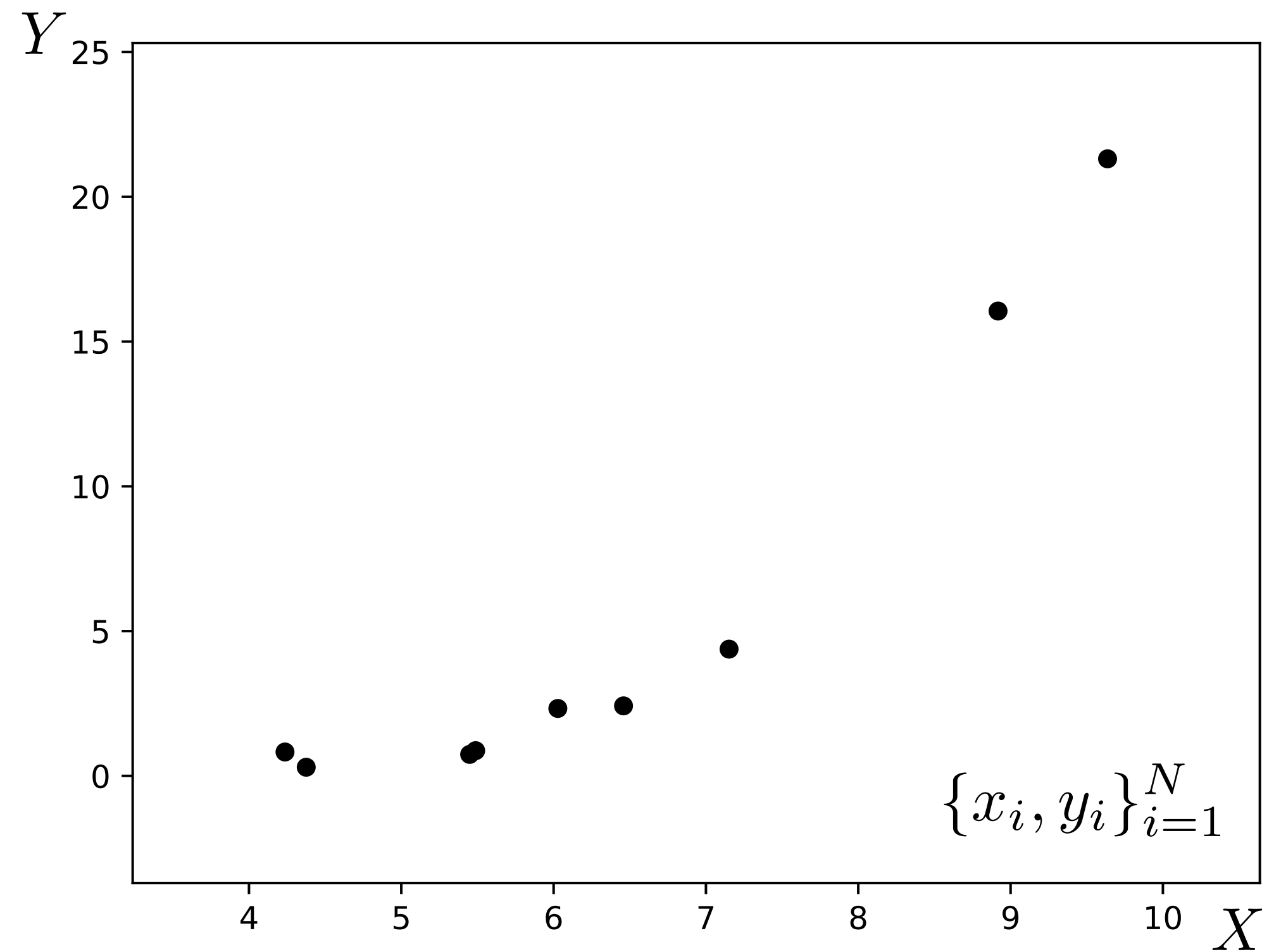
Hypothesis space

Training data

The Problem of Generalization

Linear regression

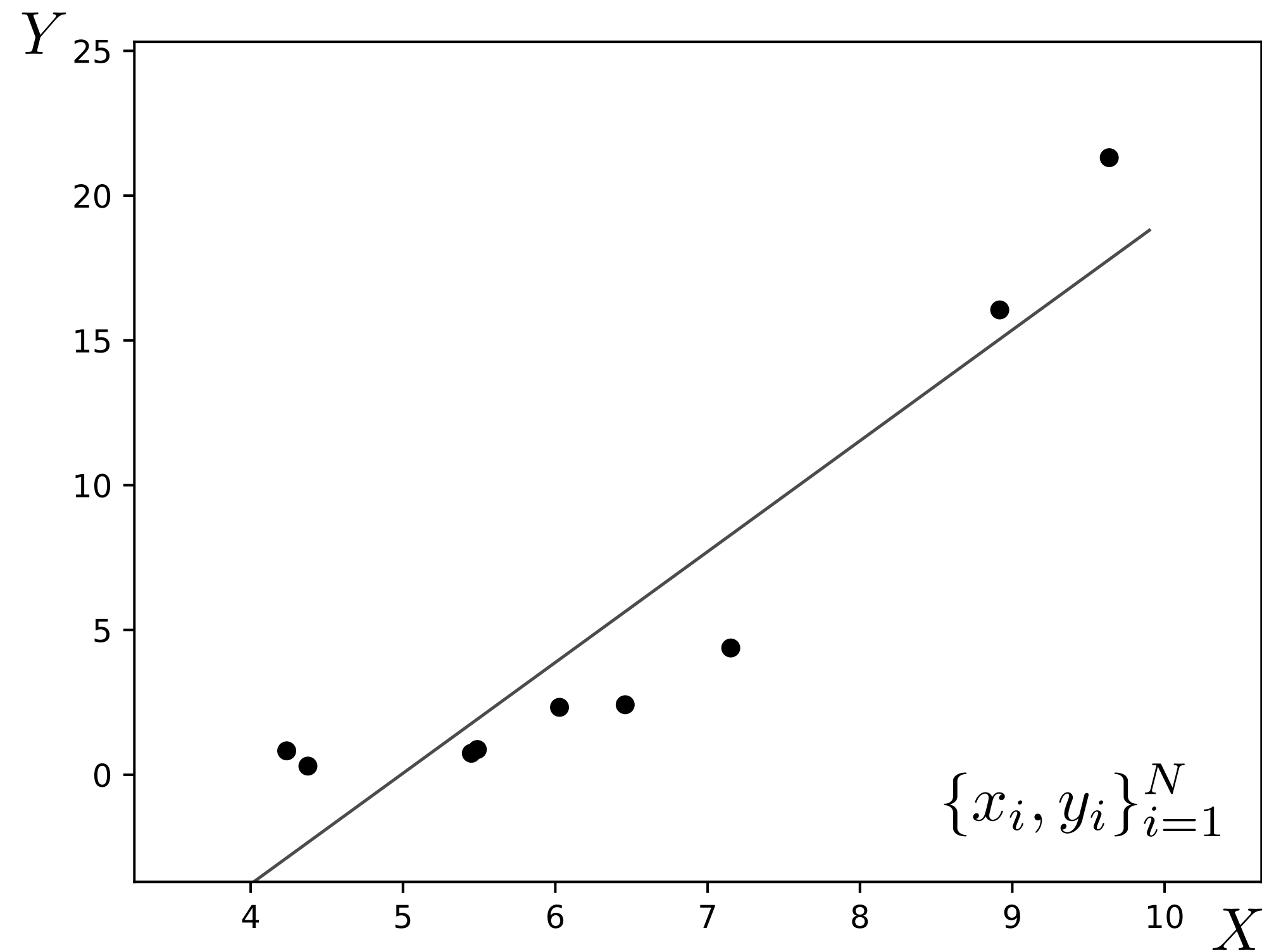
Training data



$$f_{\theta}(x) = \theta_0 + \theta_1 x$$

Linear regression

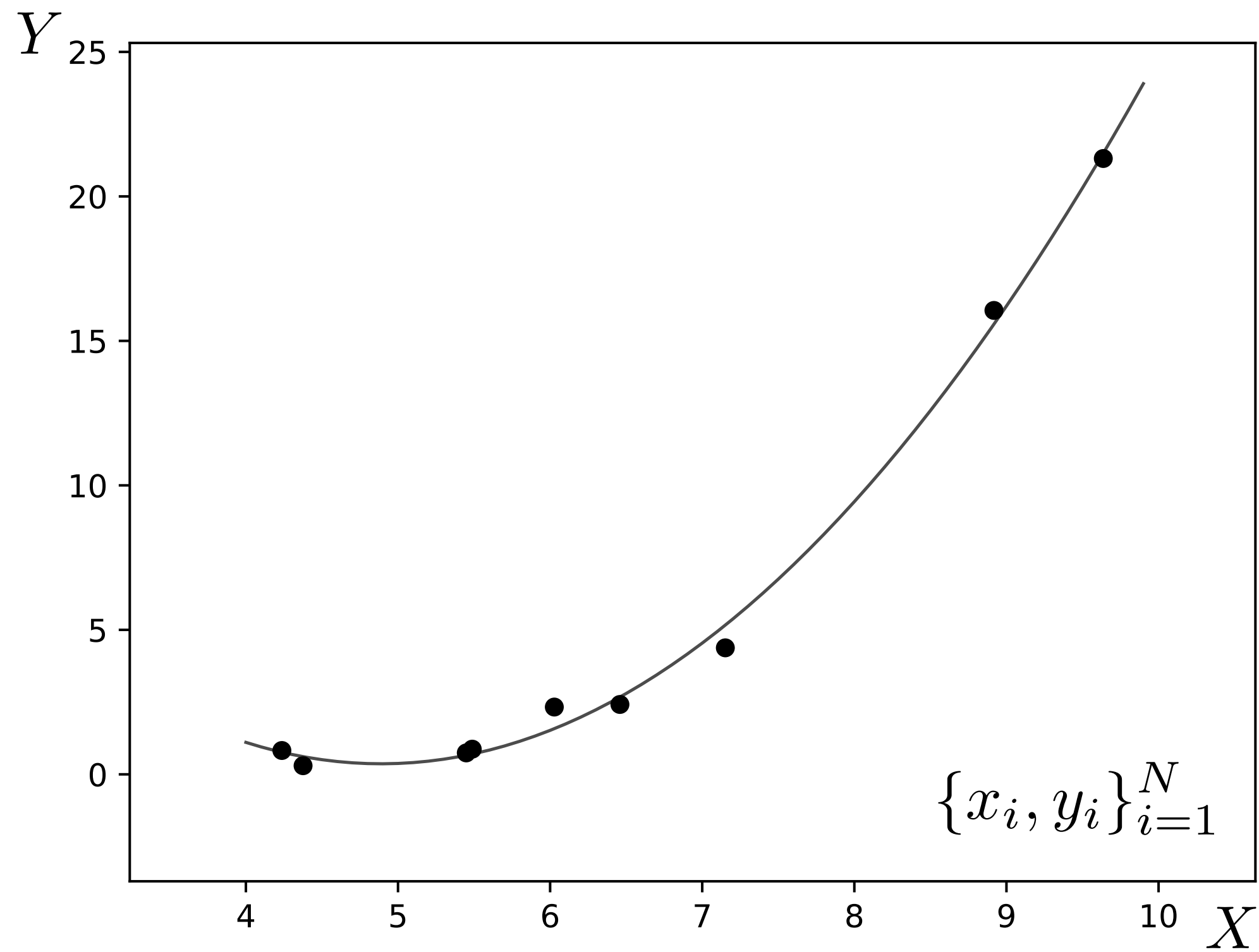
Training data



$$f_{\theta}(x) = \theta_0 + \theta_1 x$$

Polynomial regression

Training data

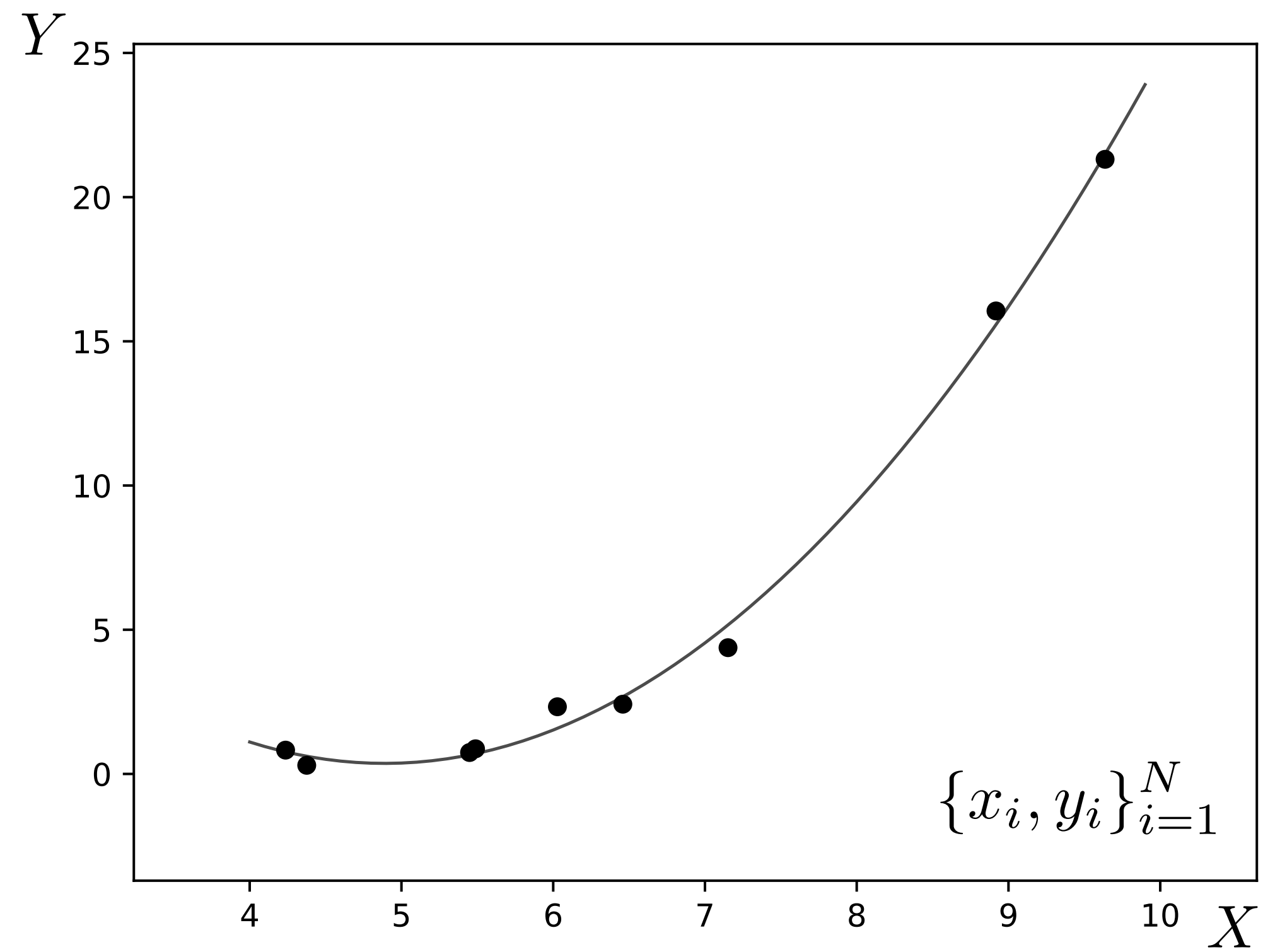


$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

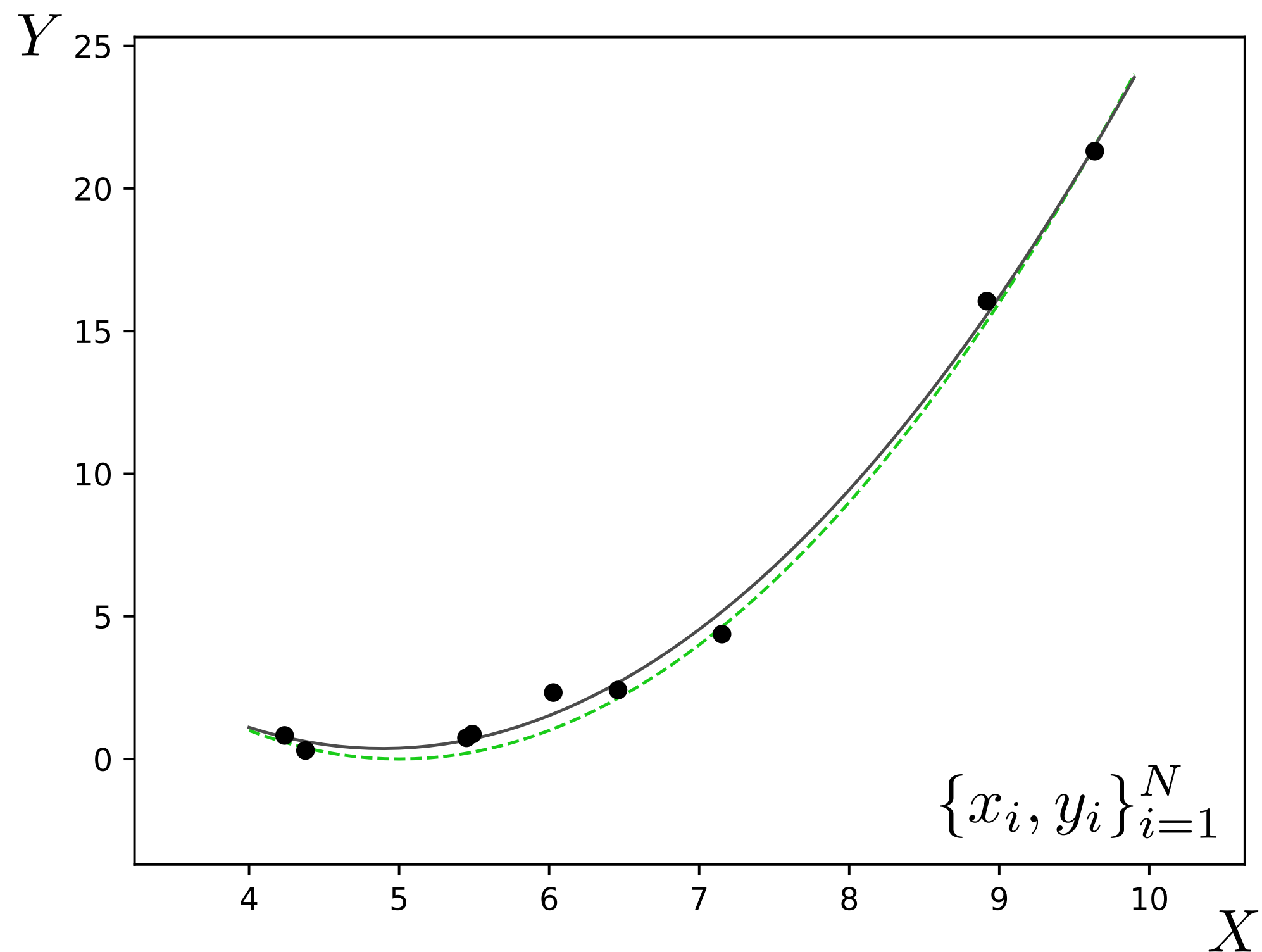
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

K-th degree polynomial regression

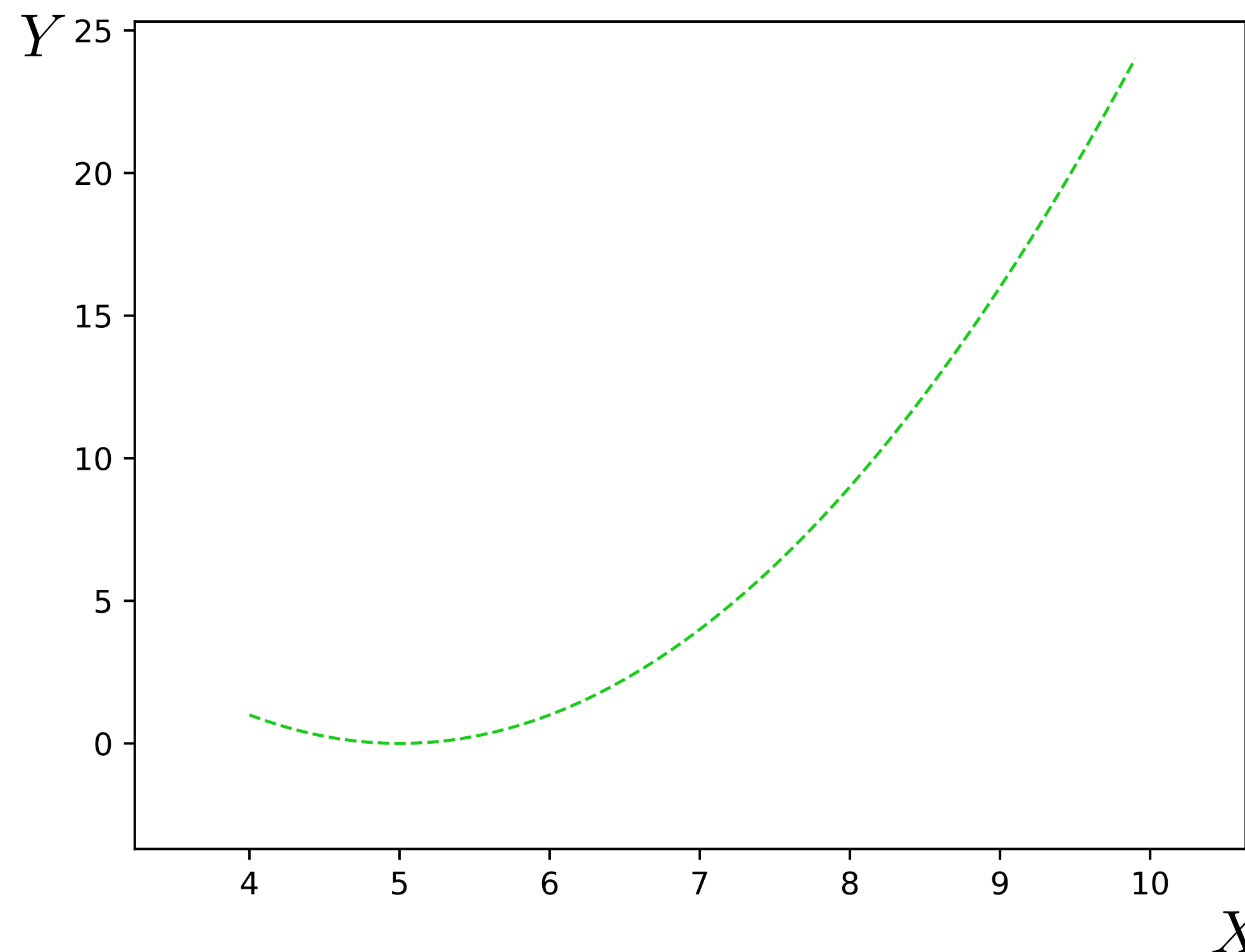
Training data



Training data



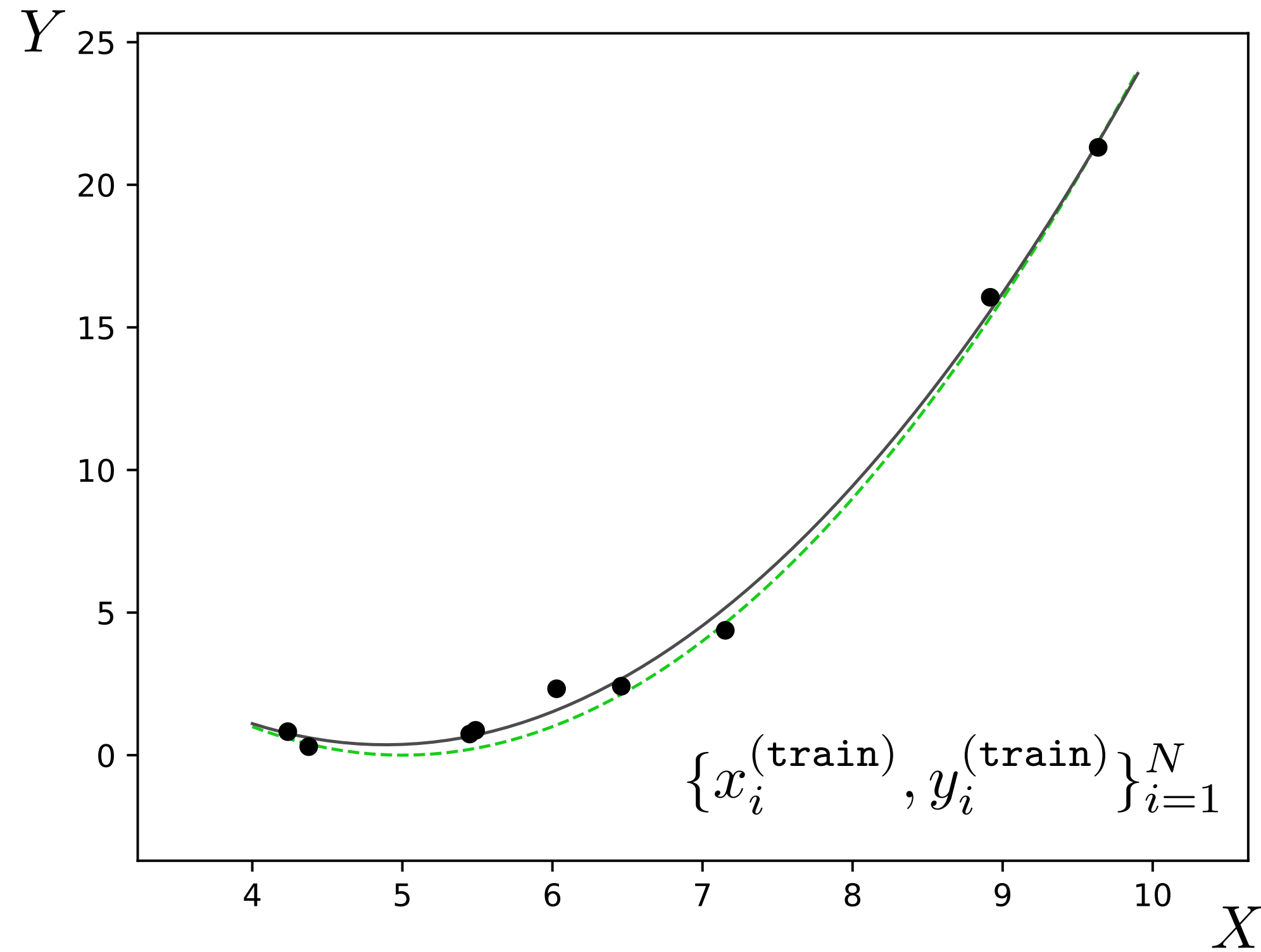
Test data



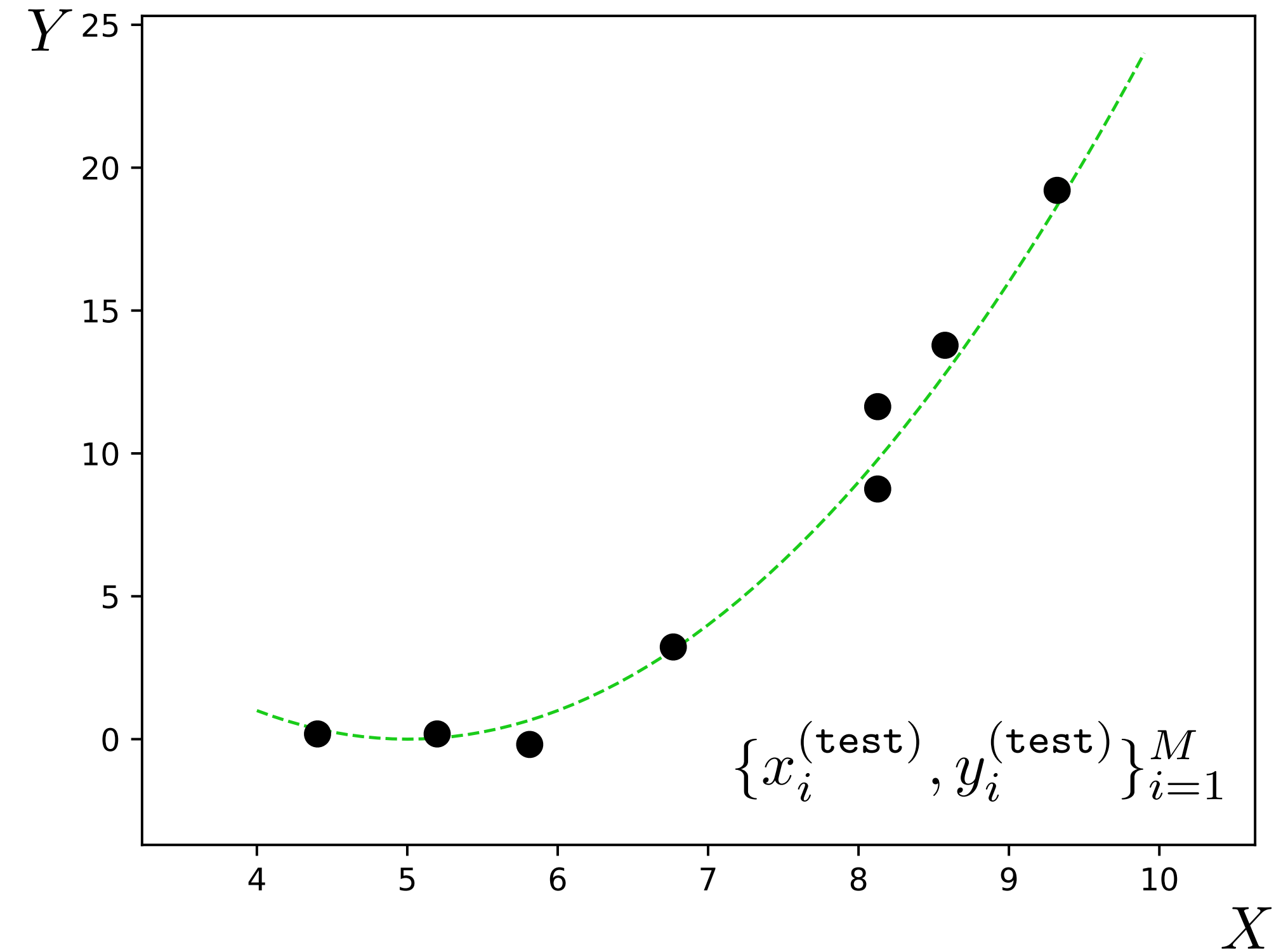
True data-generating process

p_{data}

Training data



Test data



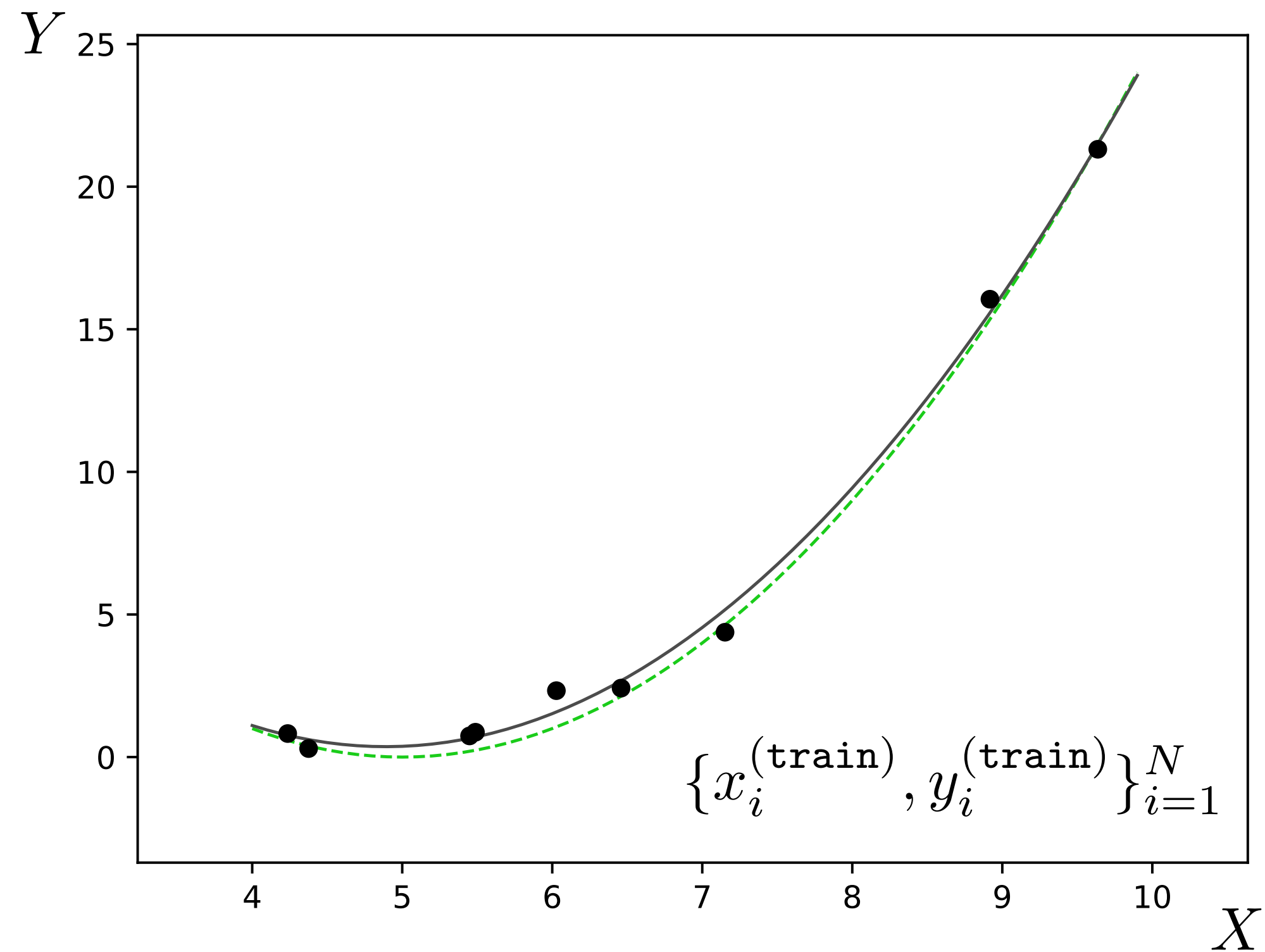
True data-generating process

p_{data}

$$\{x_i^{(\text{train})}, y_i^{(\text{train})}\}_{i=1}^N \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

$$\{x_i^{(\text{test})}, y_i^{(\text{test})}\}_{i=1}^M \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

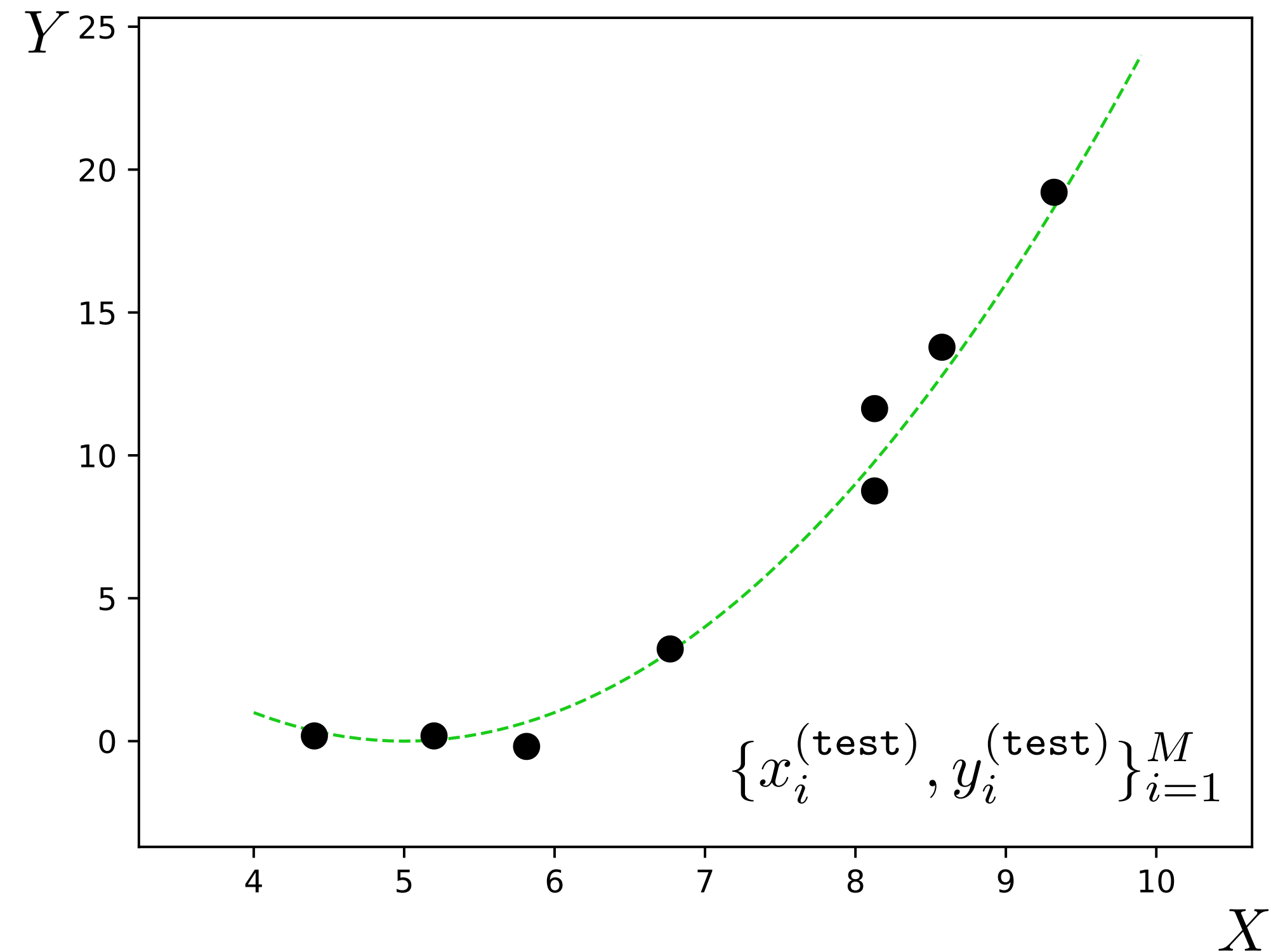
Training data



Training objective:

$$\sum_{i=1}^N (f_{\theta}(x_i^{\text{train}}) - y_i^{\text{train}})^2$$

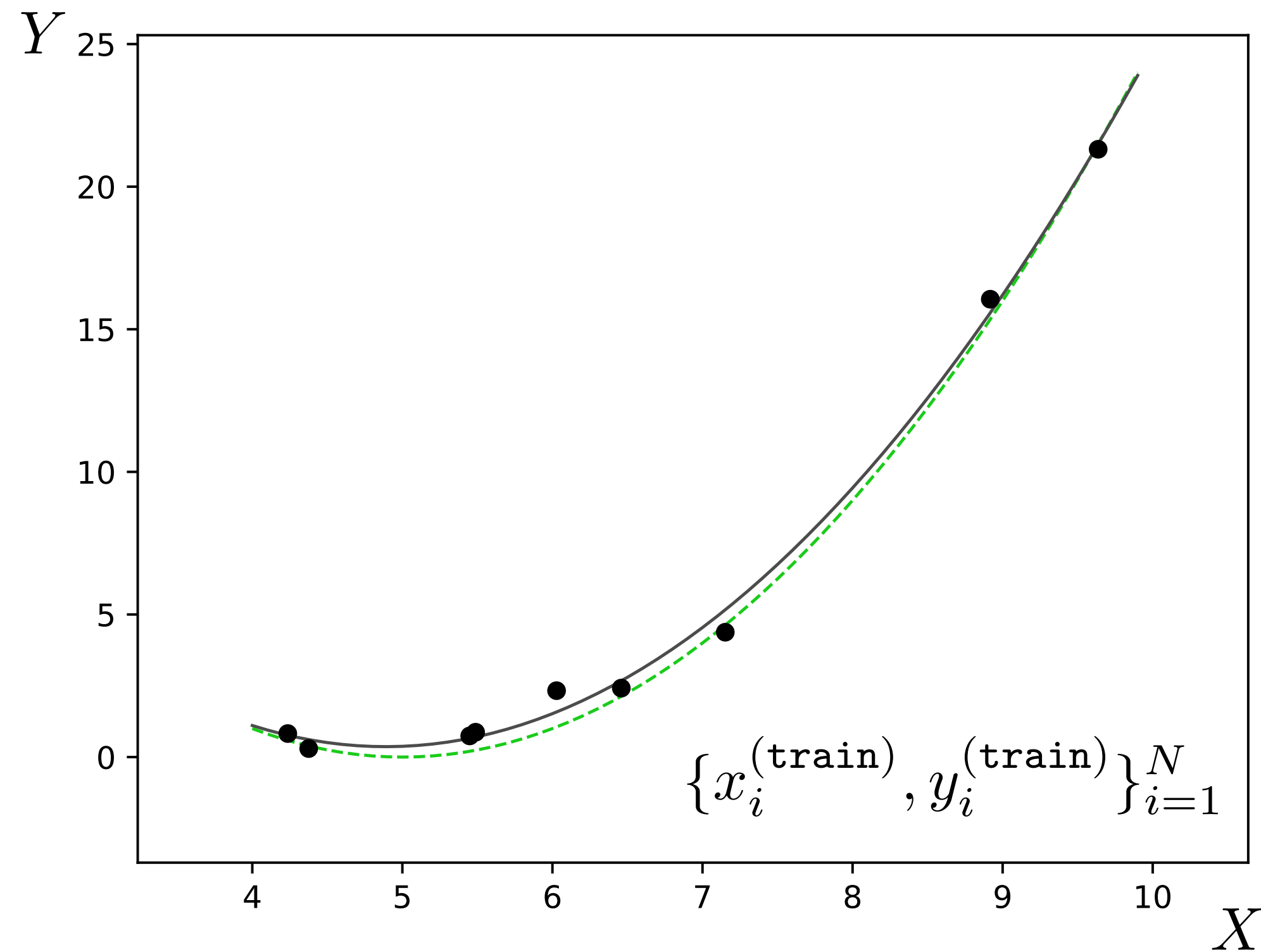
Test data



Test time evaluation:

$$\sum_{i=1}^M (f_{\theta}(x_i^{\text{test}}) - y_i^{\text{test}})^2$$

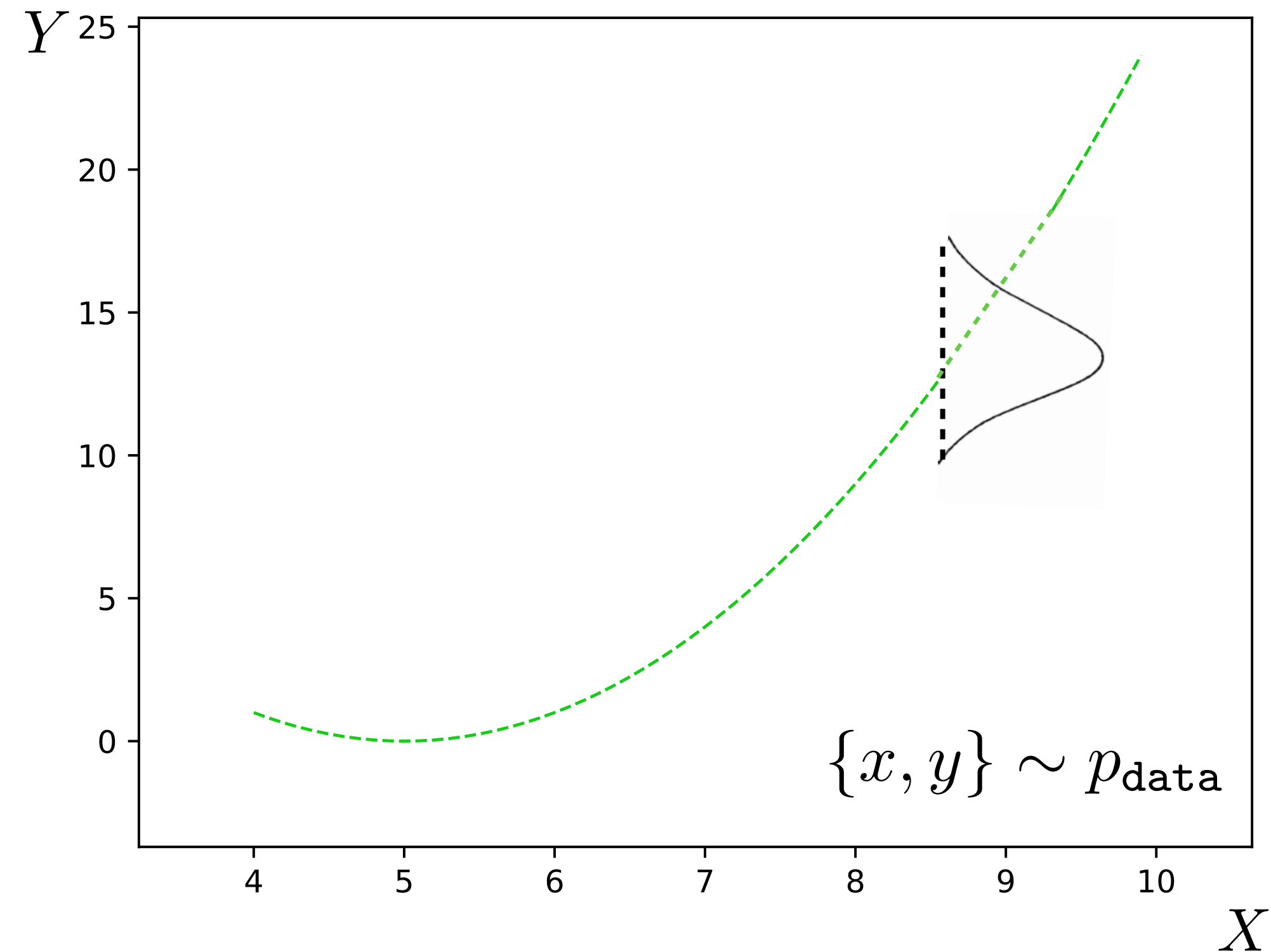
Training data



Training objective:

$$\sum_{i=1}^N (f_{\theta}(x_i^{\text{train}}) - y_i^{\text{train}})^2$$

Test data

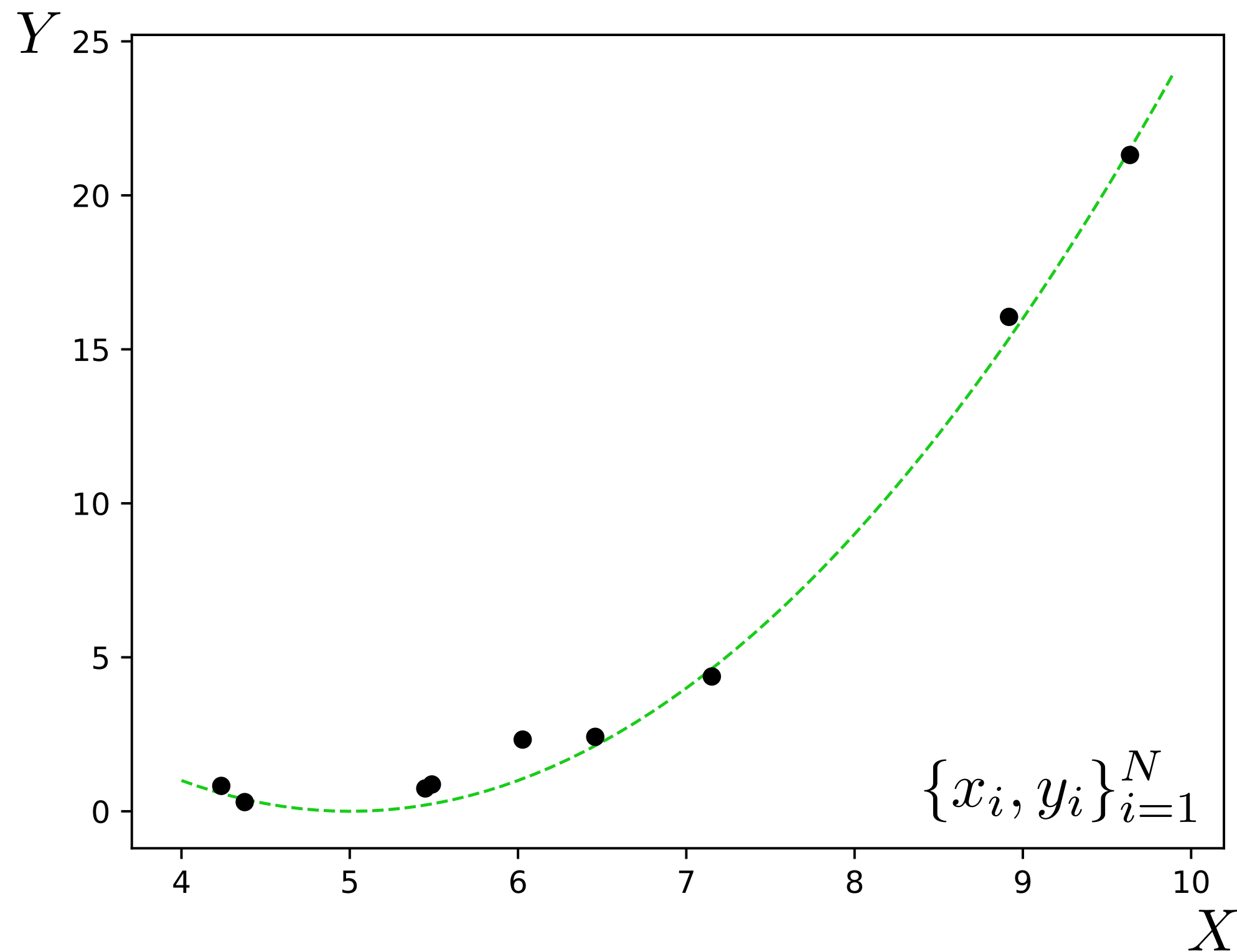


True objective:

$$\mathbb{E}_{\{x, y\} \sim p_{\text{data}}} [(f_{\theta}(x) - y)^2]$$

What happens as we add more basis functions?

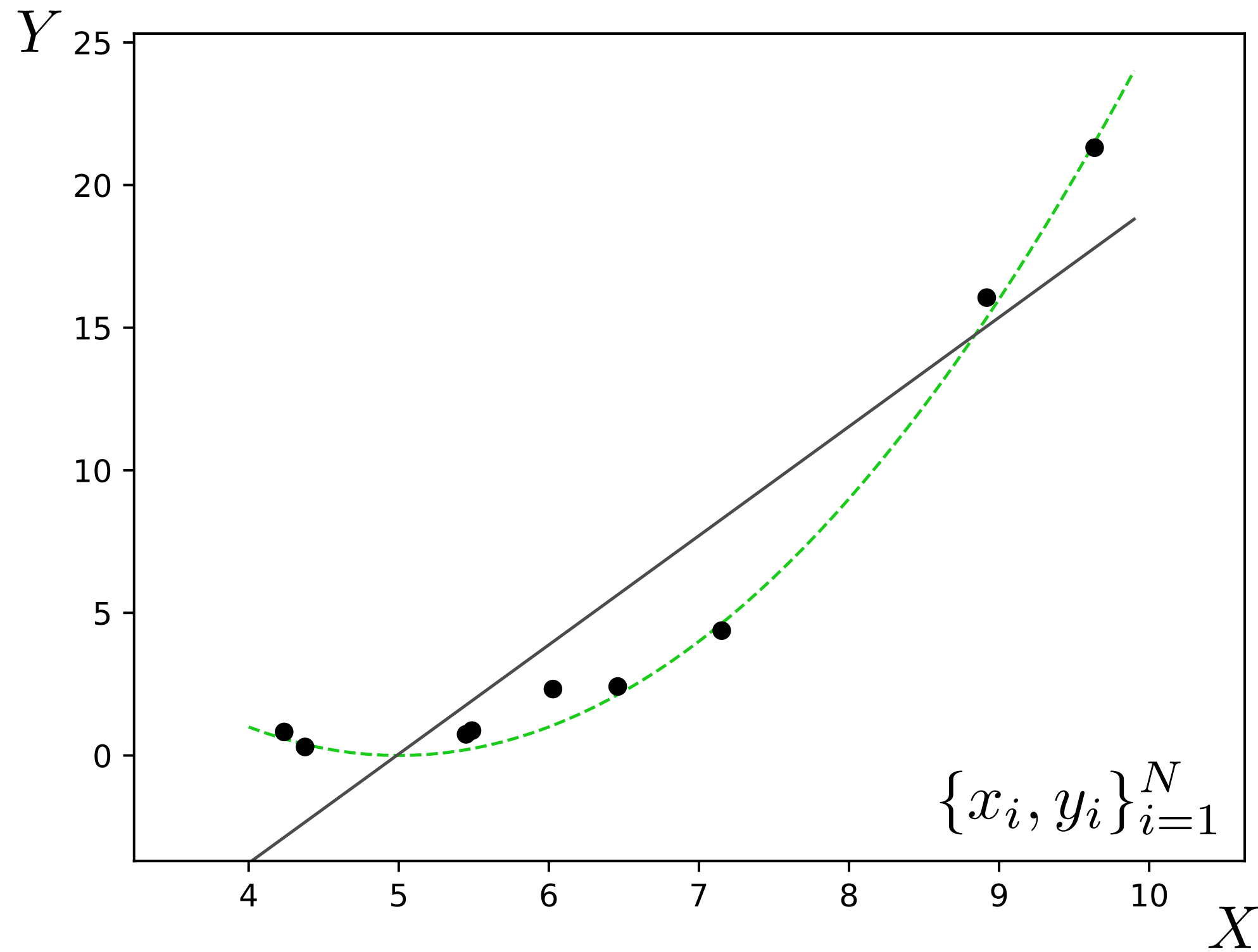
Training data



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

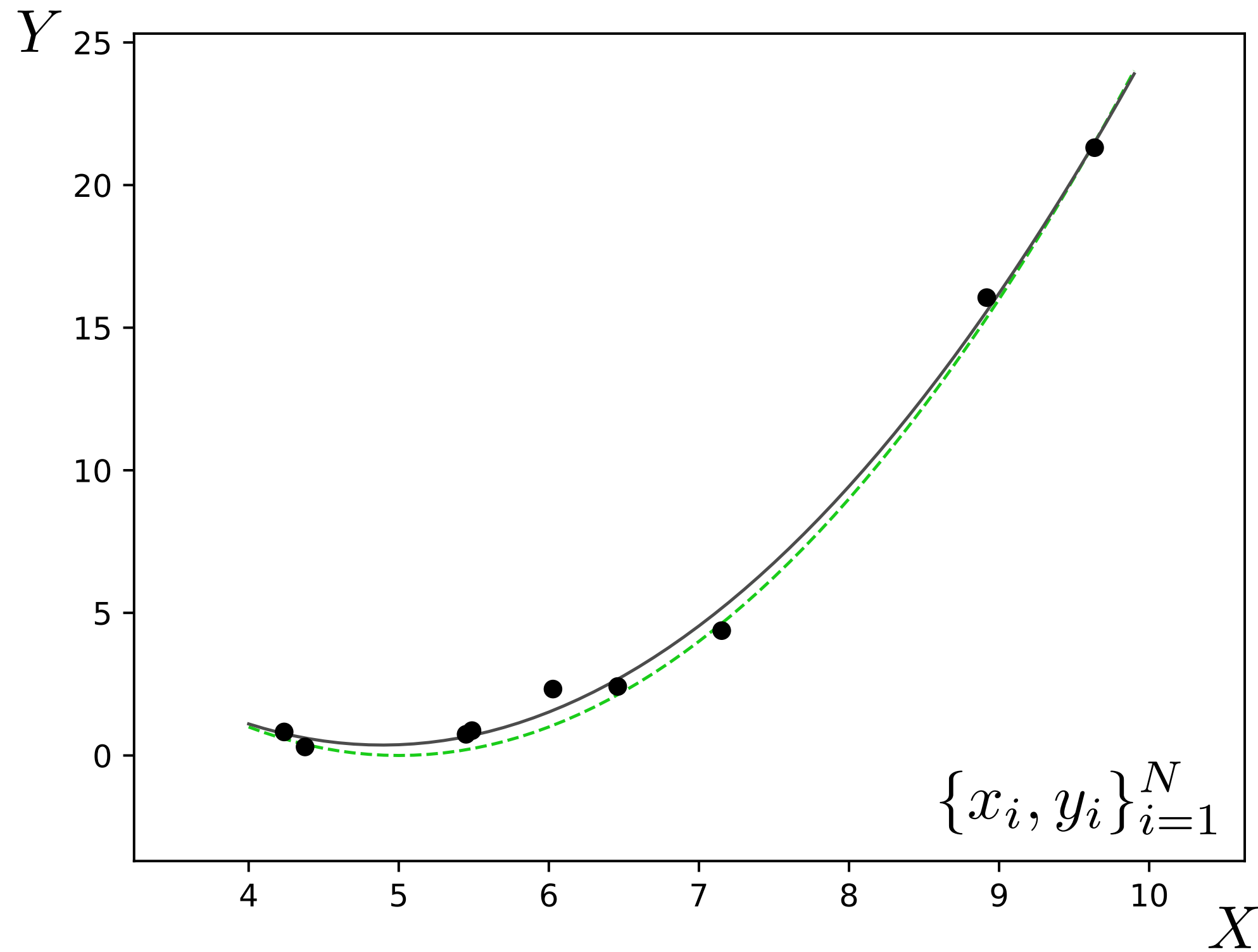
K = 1



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

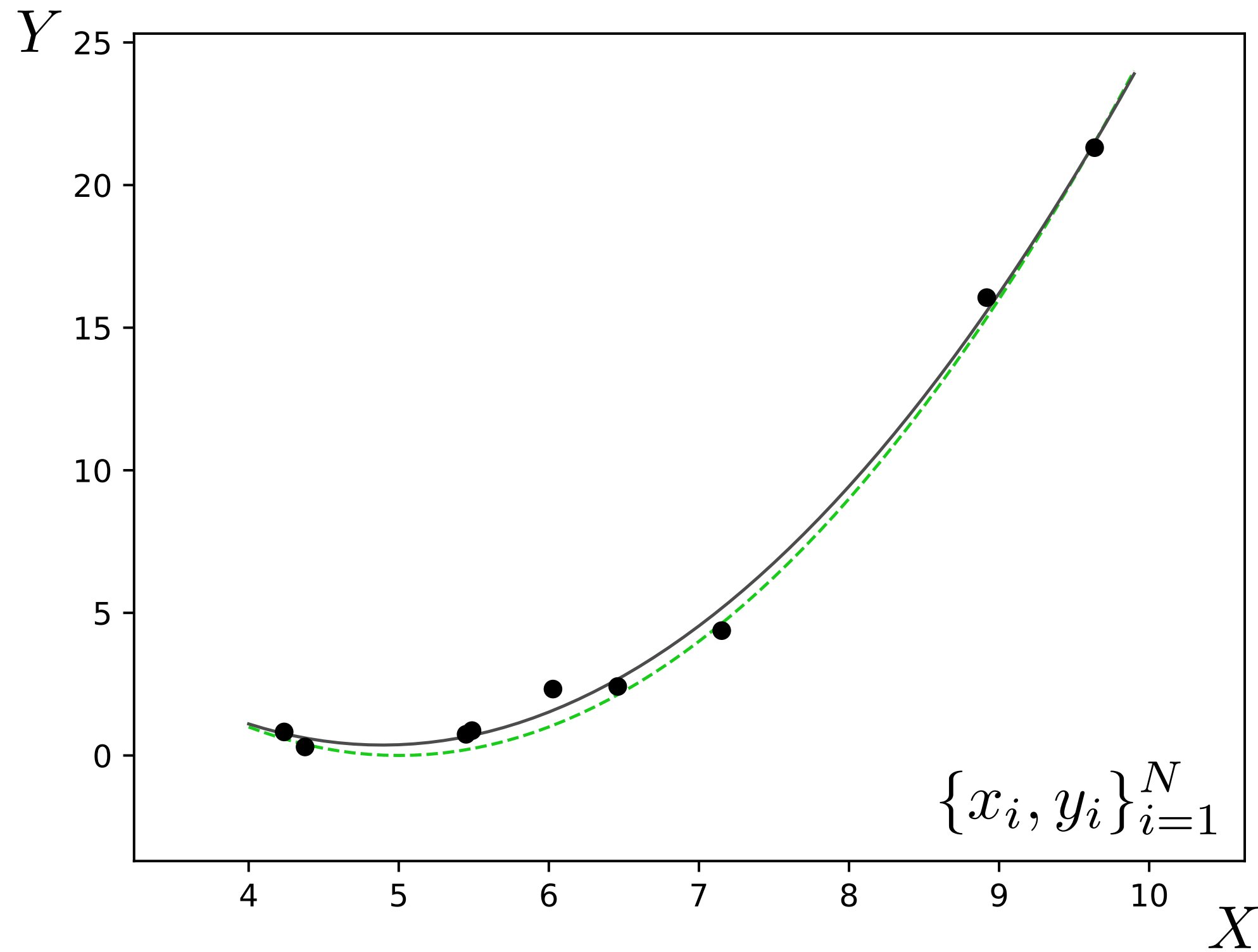
K = 2



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

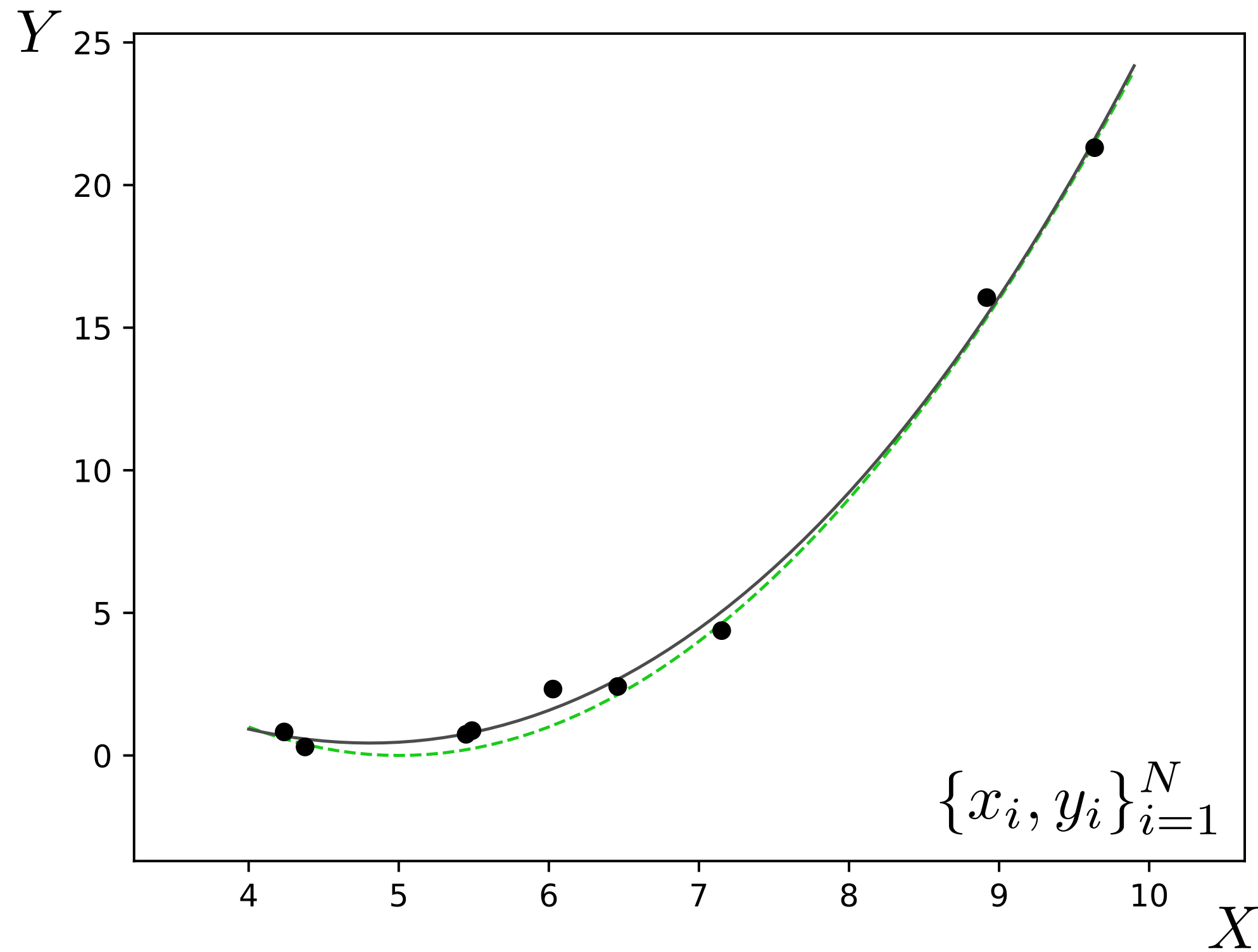
K = 3



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

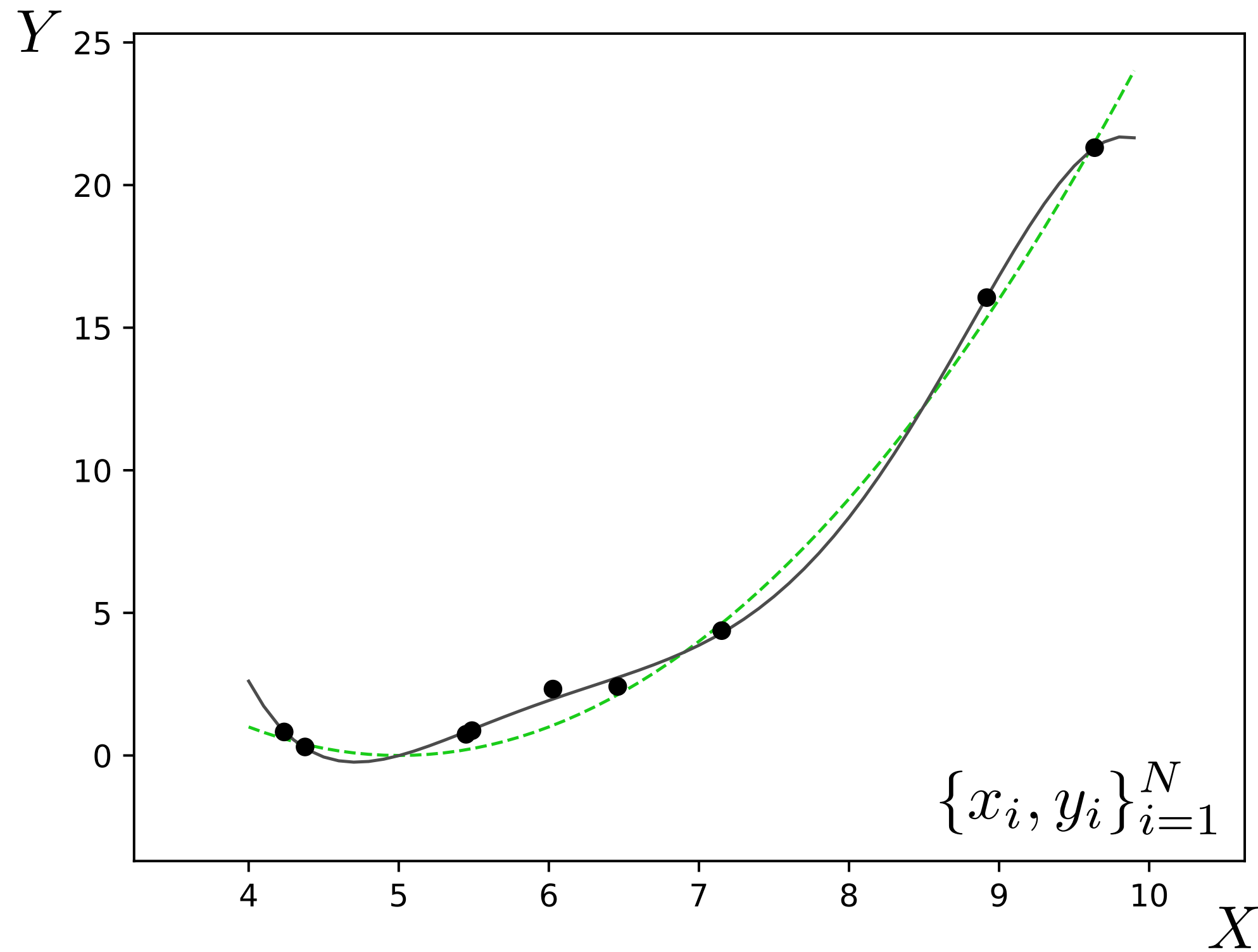
K = 4



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

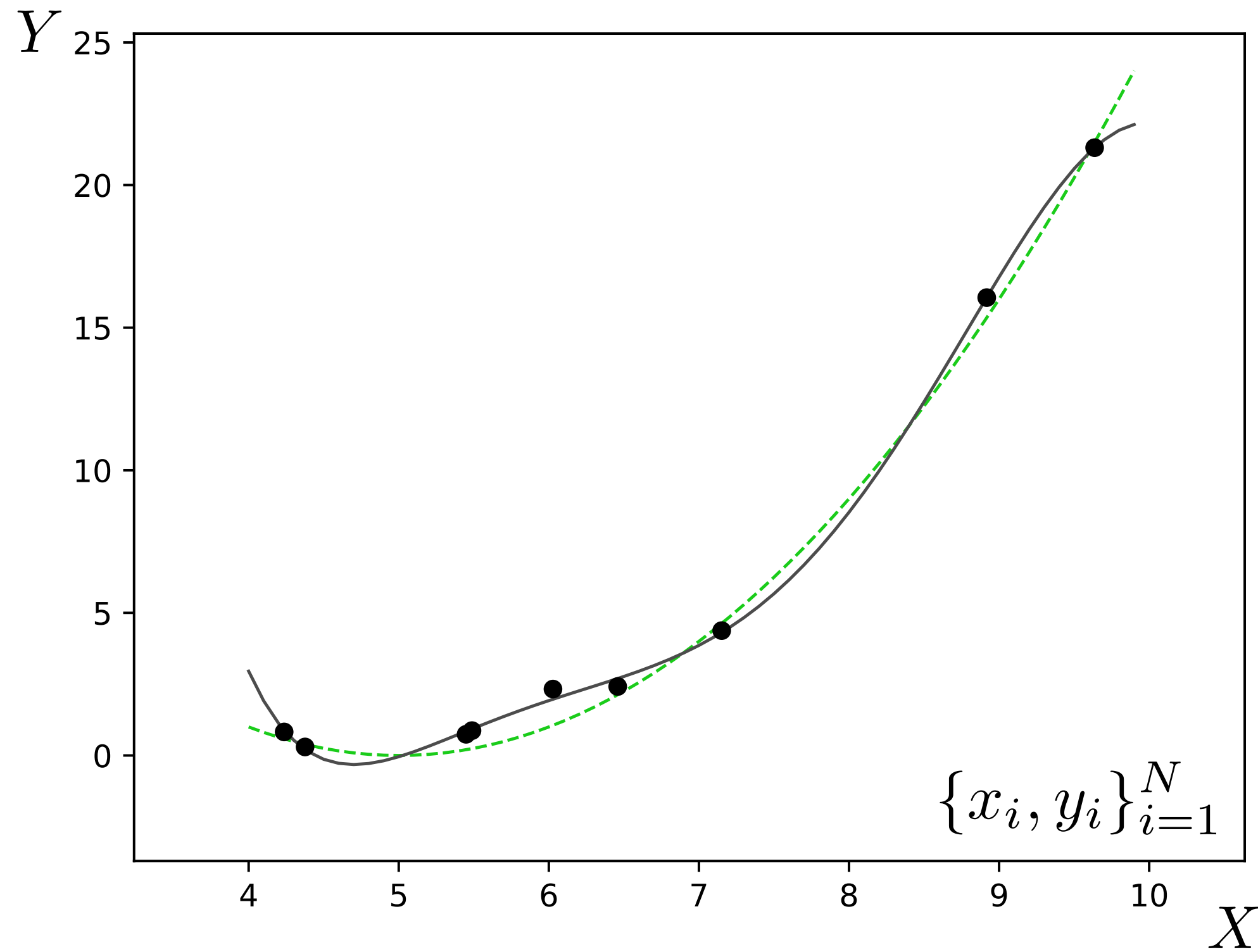
K = 5



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

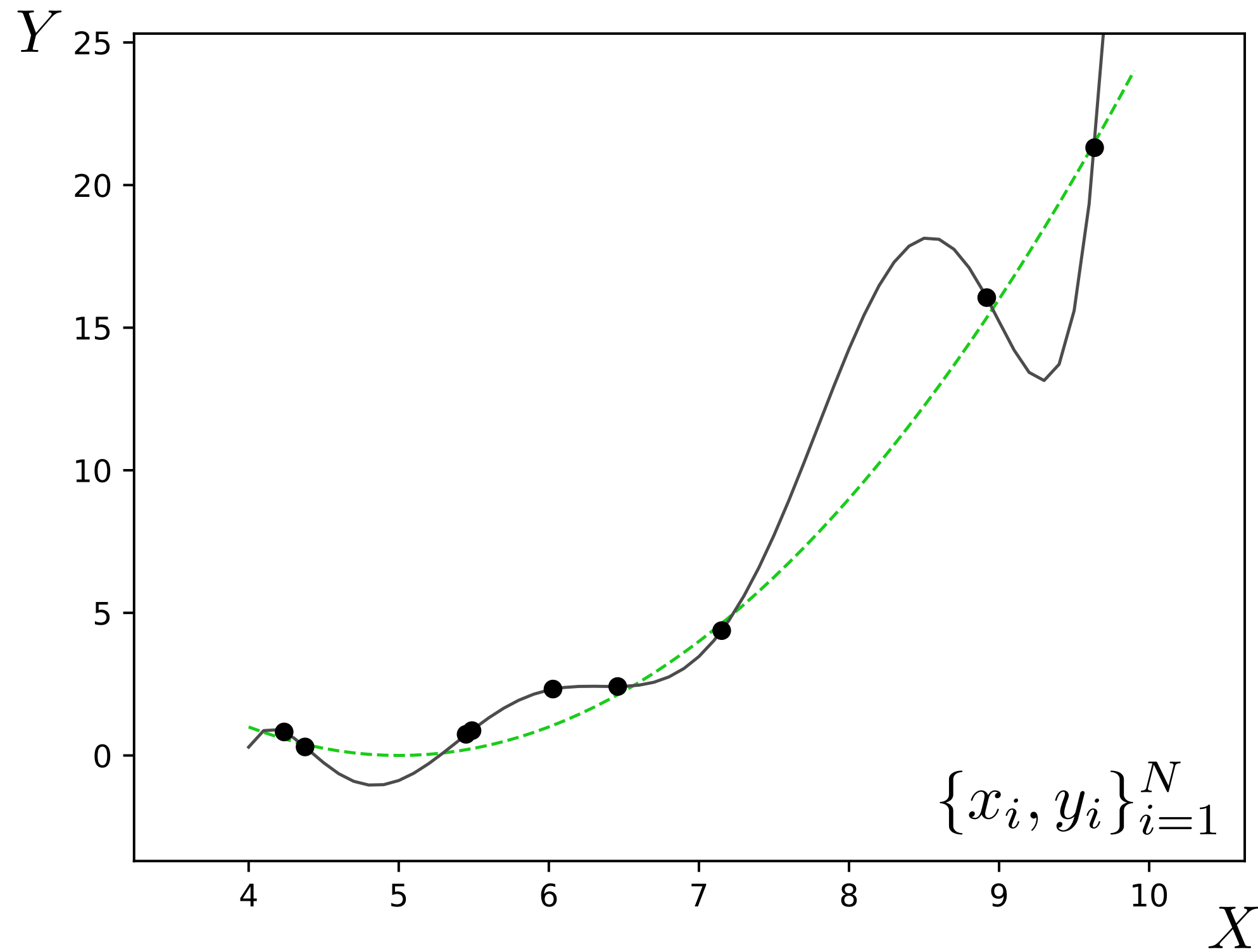
K = 6



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

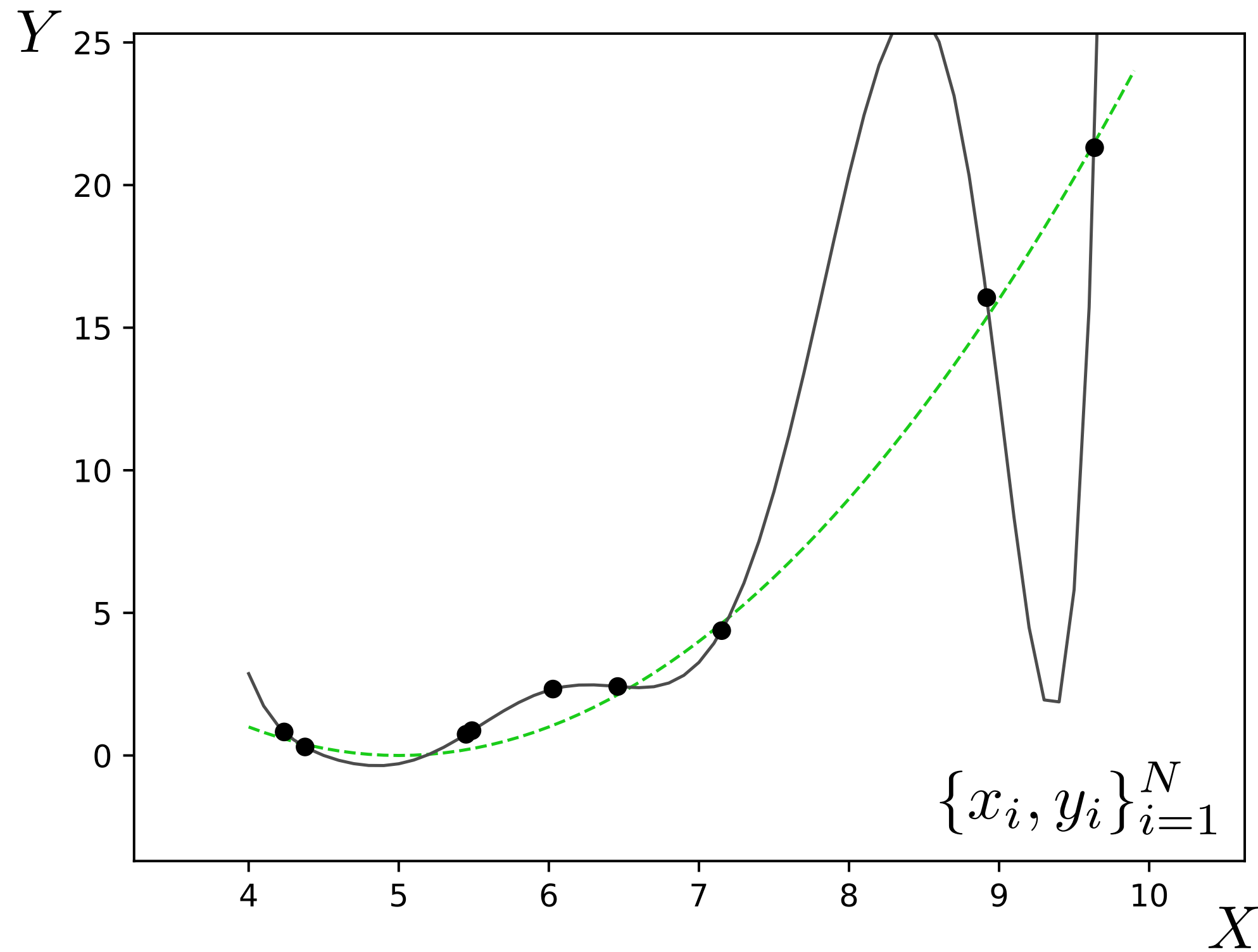
K = 7



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

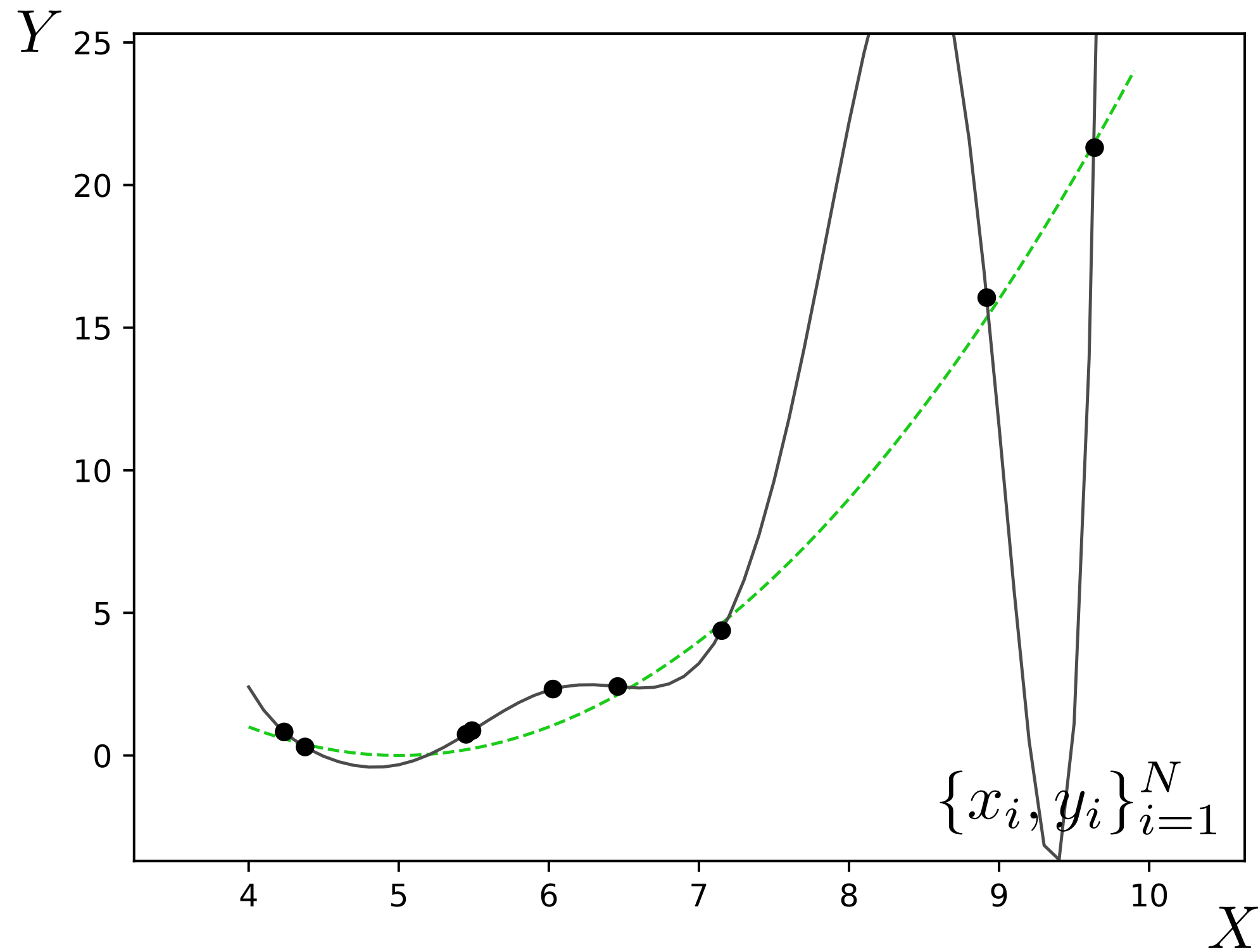
K = 8



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

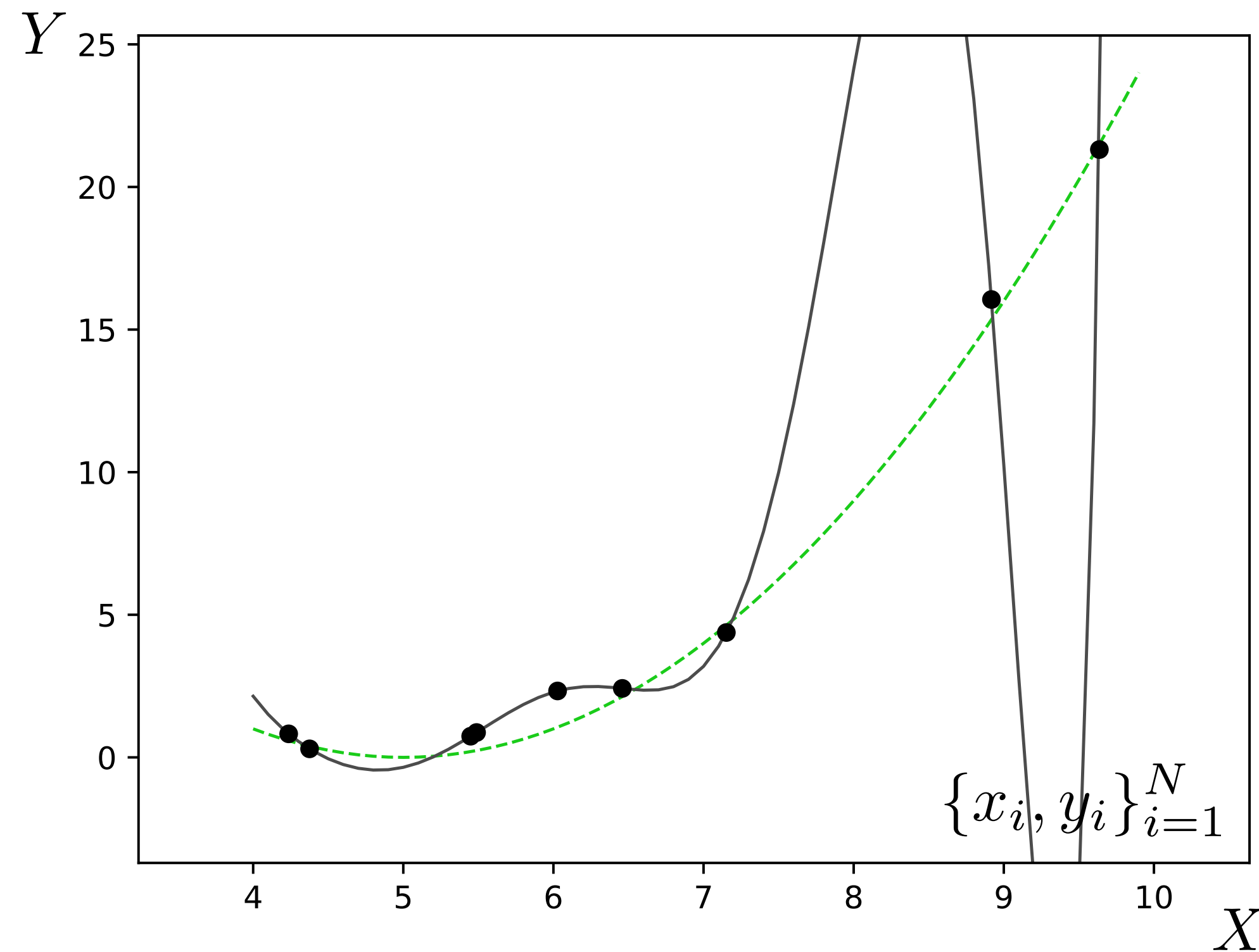
K = 9



$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?

K = 10

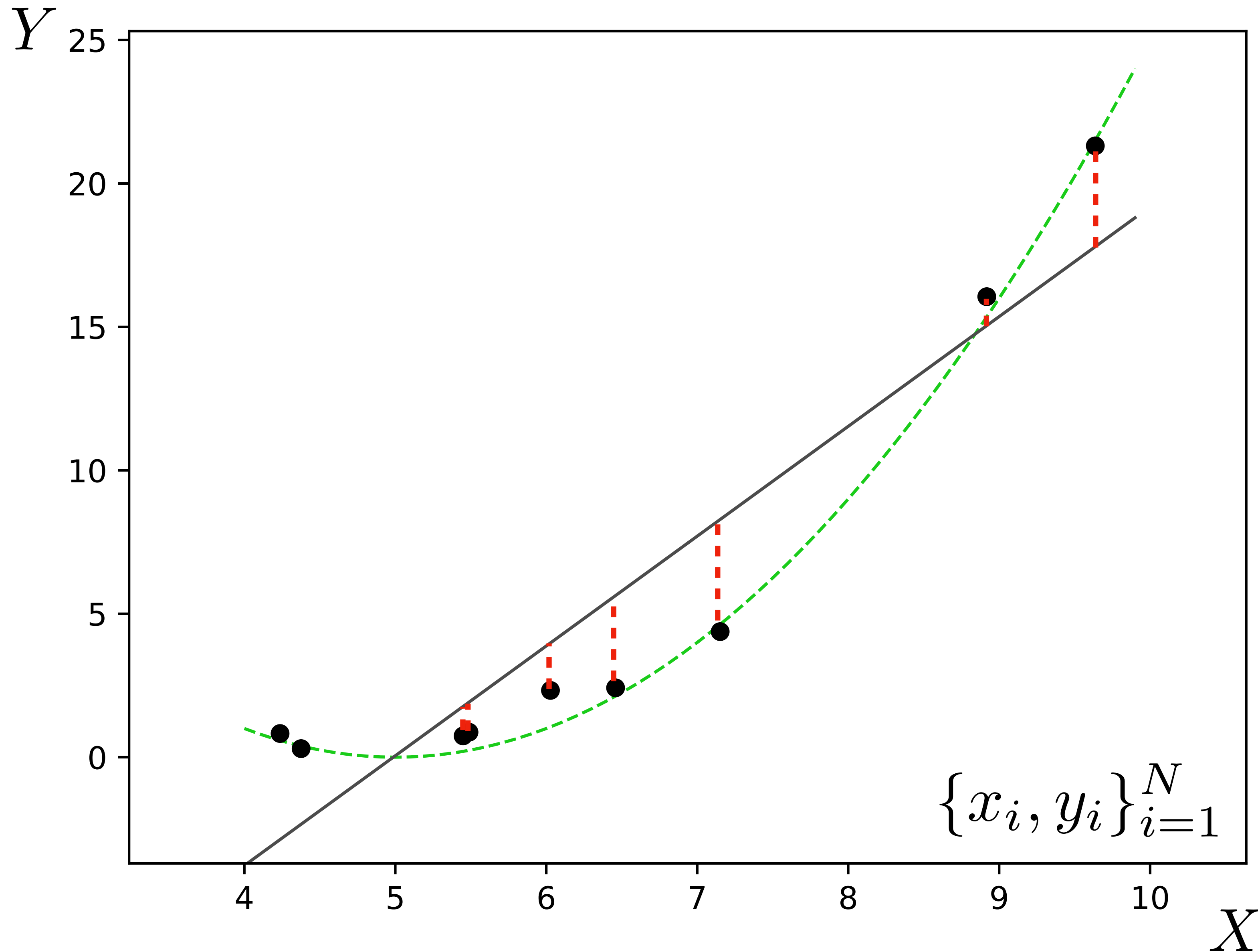


$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

This phenomenon is called **overfitting**.

It occurs when we have too high **capacity** a model, e.g., too many free parameters, too few data points to pin these parameters down.

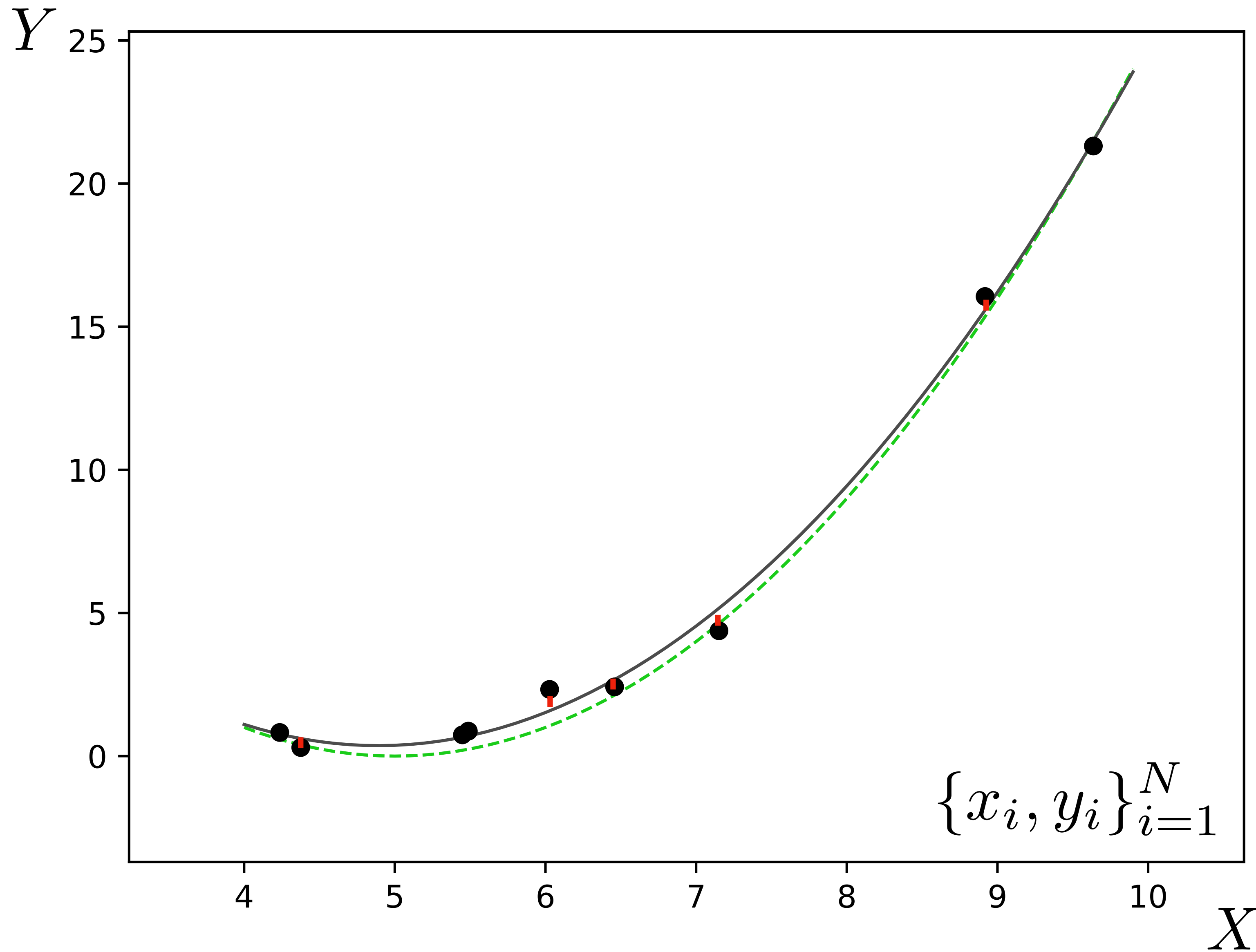
$$K = 1$$



When the model does not have the capacity to capture the true function, we call this **underfitting**.

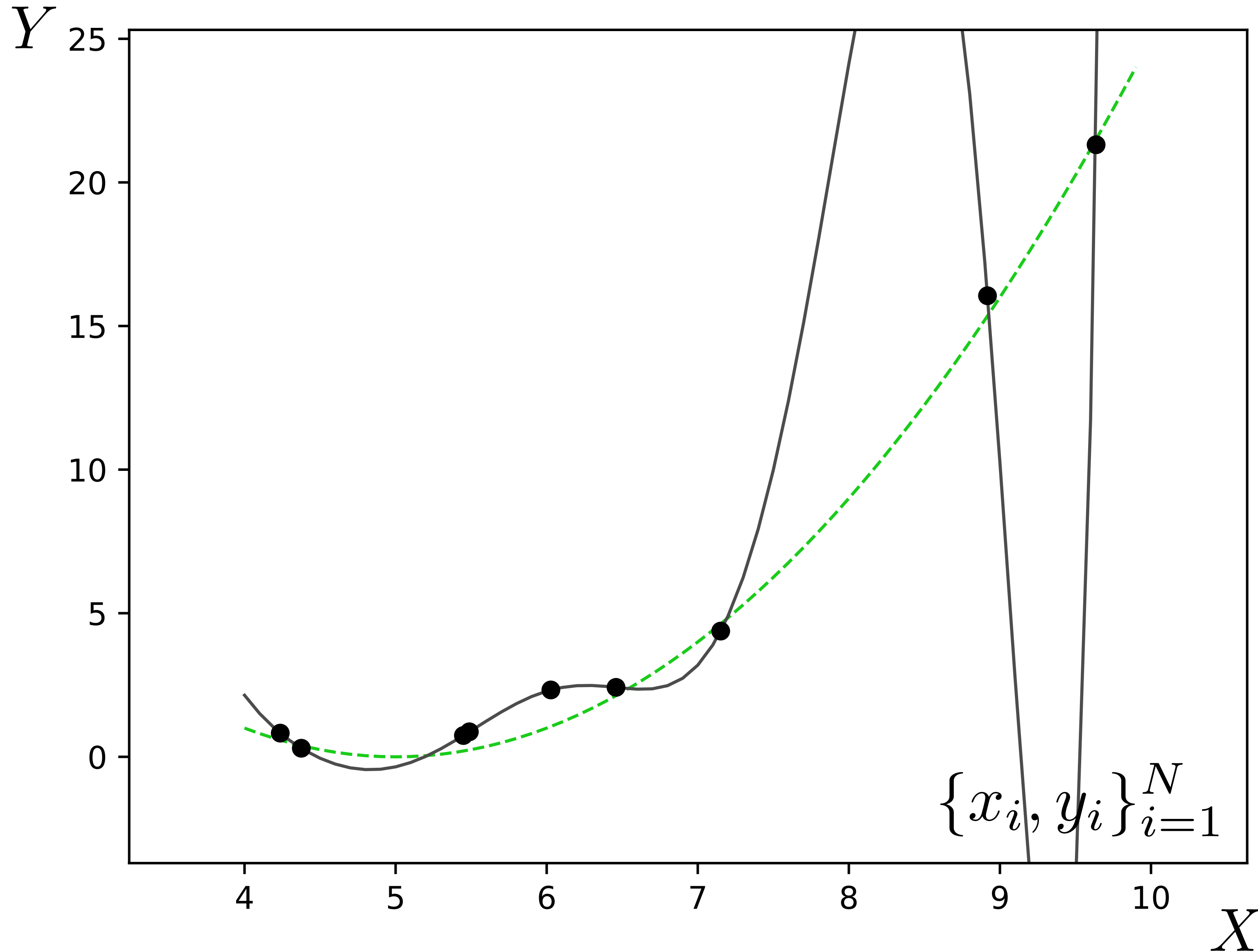
An underfit model will have high **error** on the training points. This error is known as **approximation error**.

$K = 2$



The true function is a quadratic, so a quadratic model ($K=2$) fits well.

$K = 10$

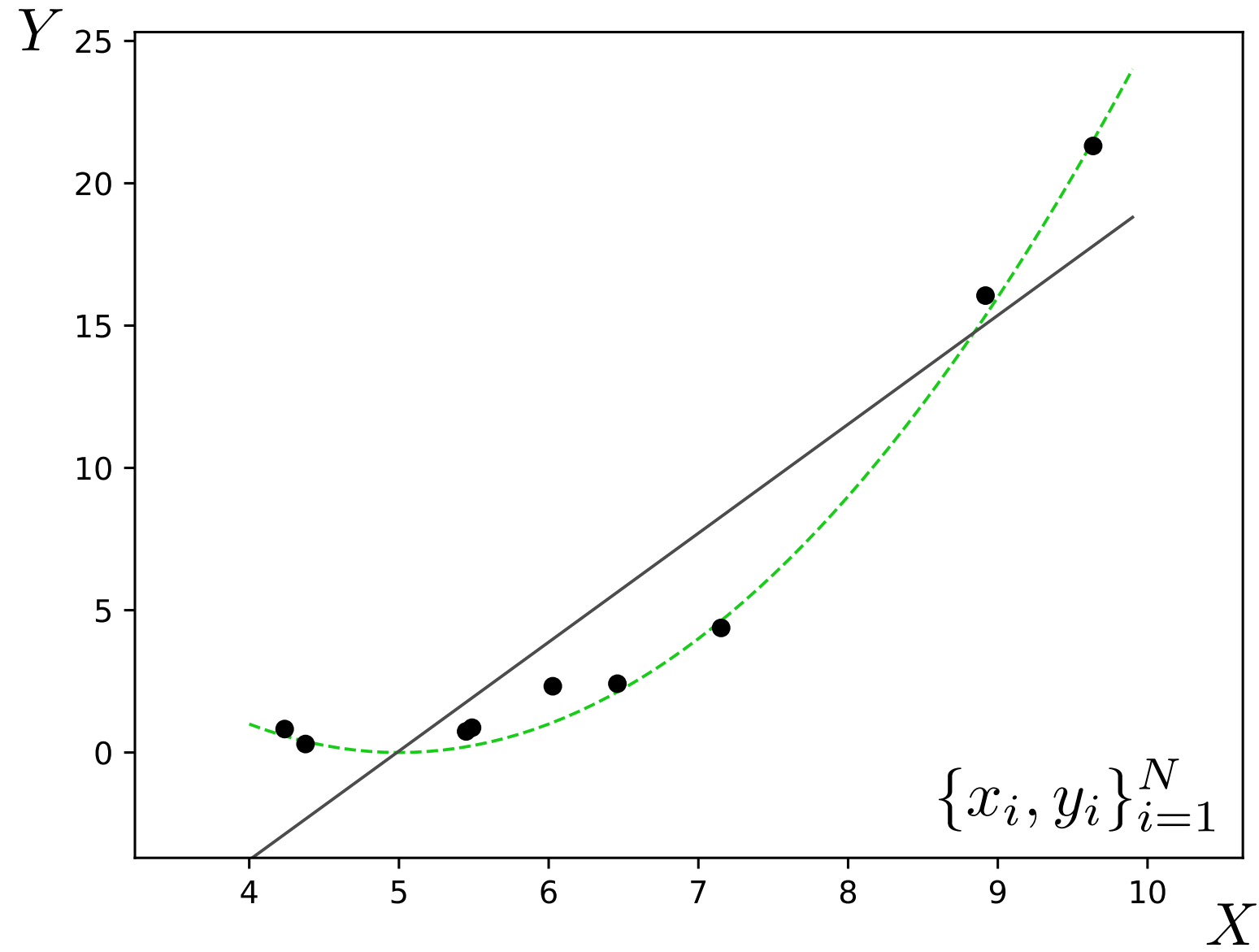


Now we have zero approximation error — the curve passes exactly through each training point.

But we have high **generalization error**, reflected in the gap between the true function and the fit line. We want to do well on *novel* queries, which will be sampled from the green curve (plus noise).

Underfitting

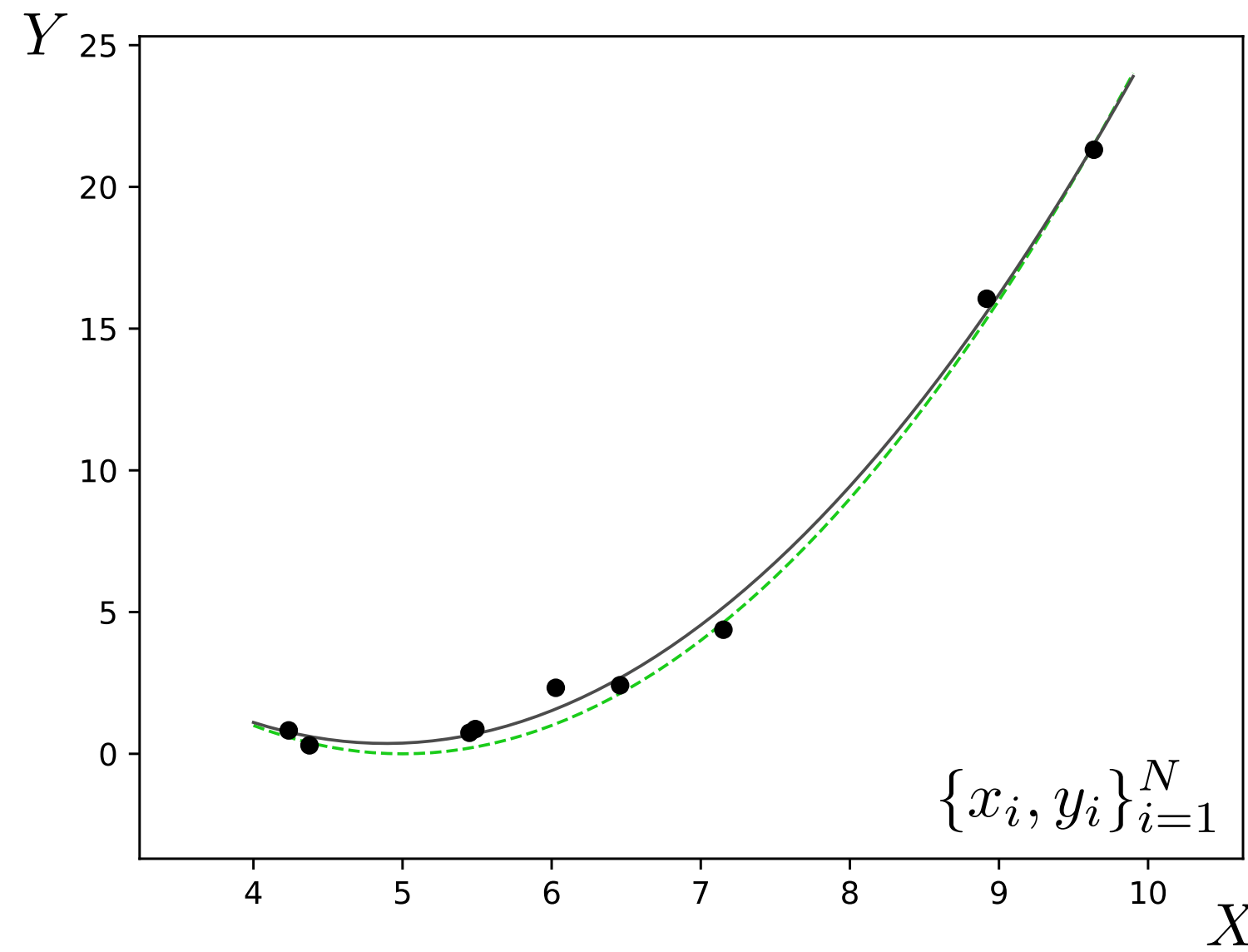
$$K = 1$$



High error on train set
High error on test set

Appropriate model

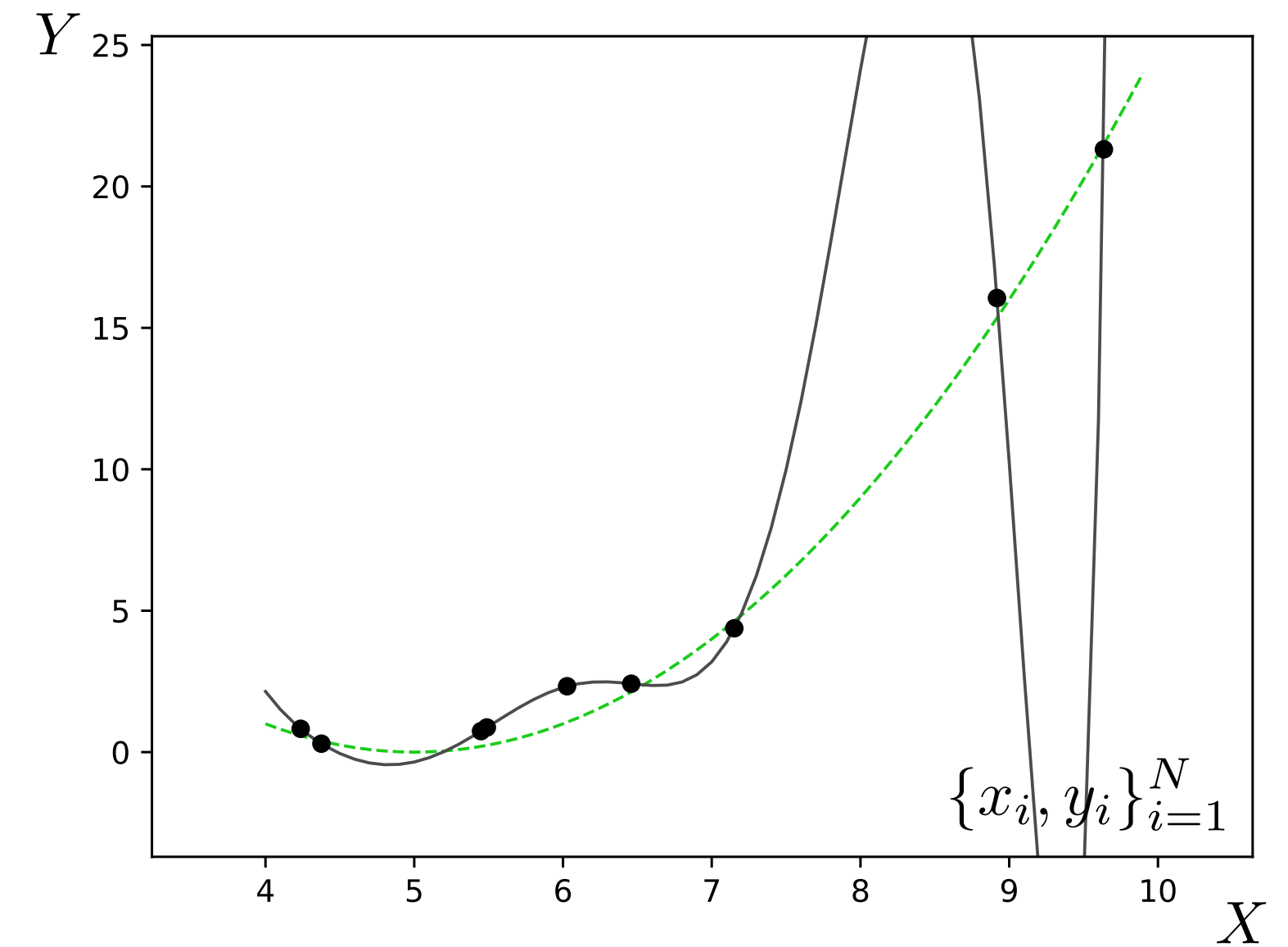
$$K = 2$$



Low error on train set
Low error on test set

Overfitting

$$K = 10$$



Lowest error on train set
High error on test set

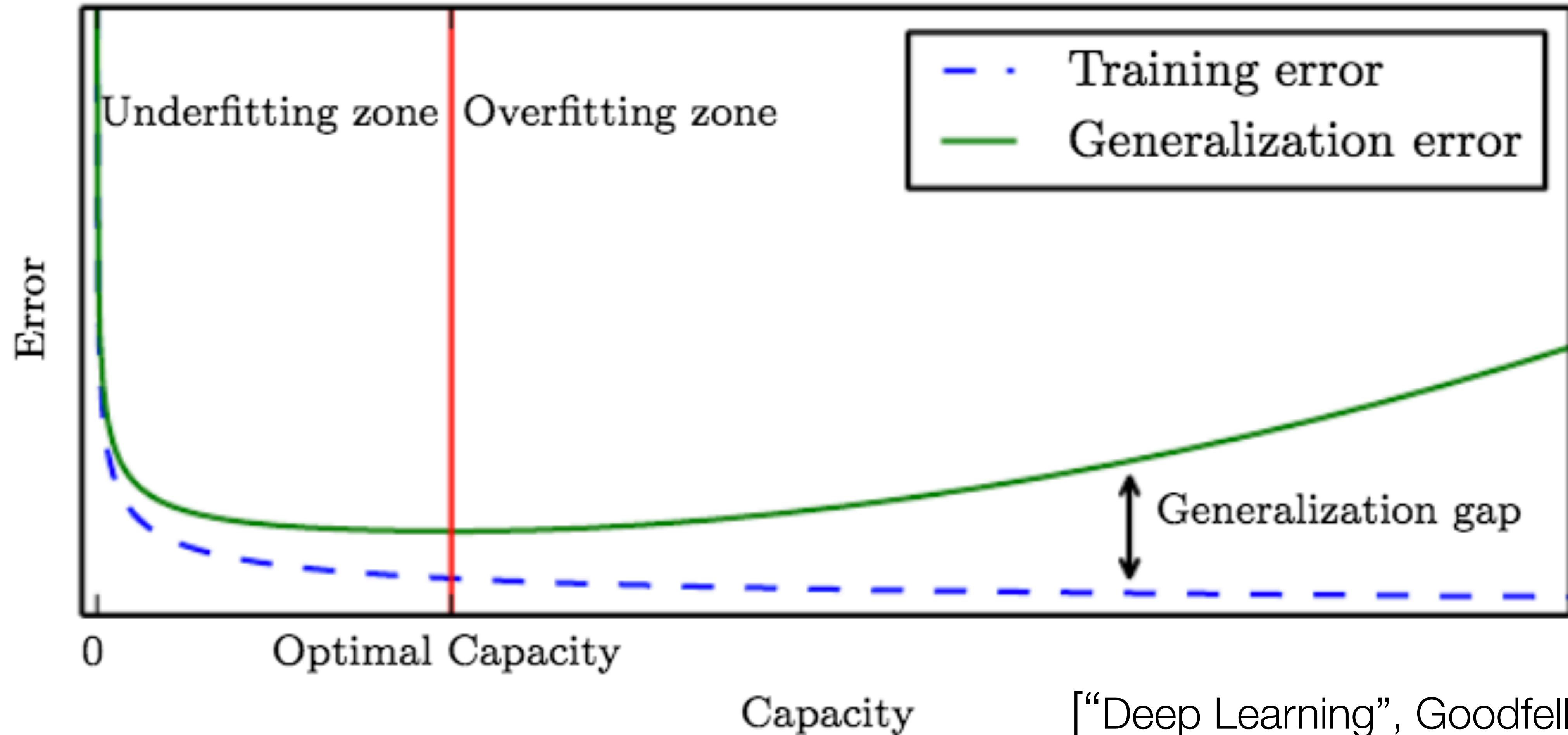
We need to control the **capacity** of the model (e.g., use the appropriate number of free parameters).

The capacity may be defined as the number of hypotheses under consideration in the hypothesis space.

Complex models with many free parameters have high capacity.

Simple models have low capacity.

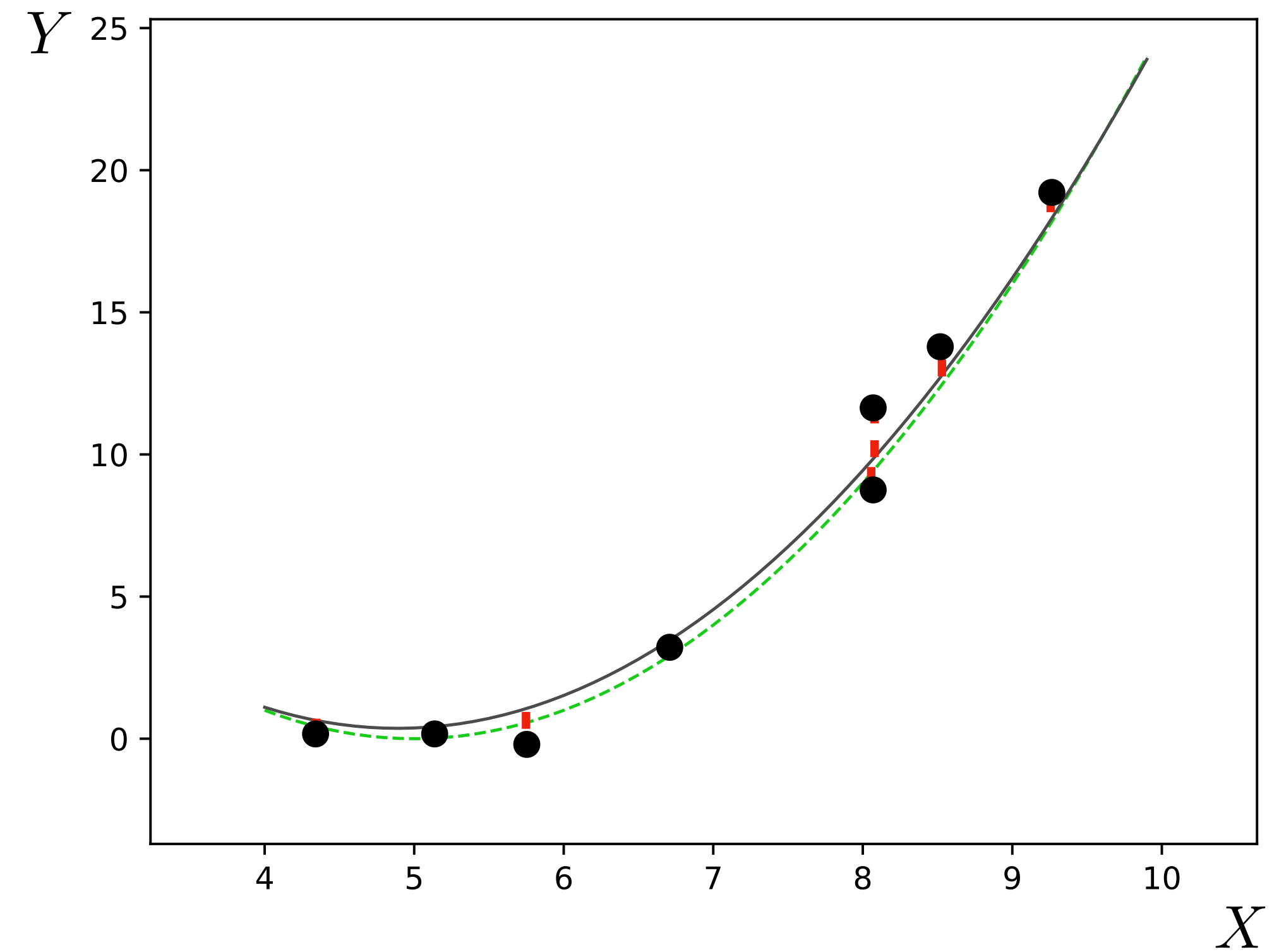
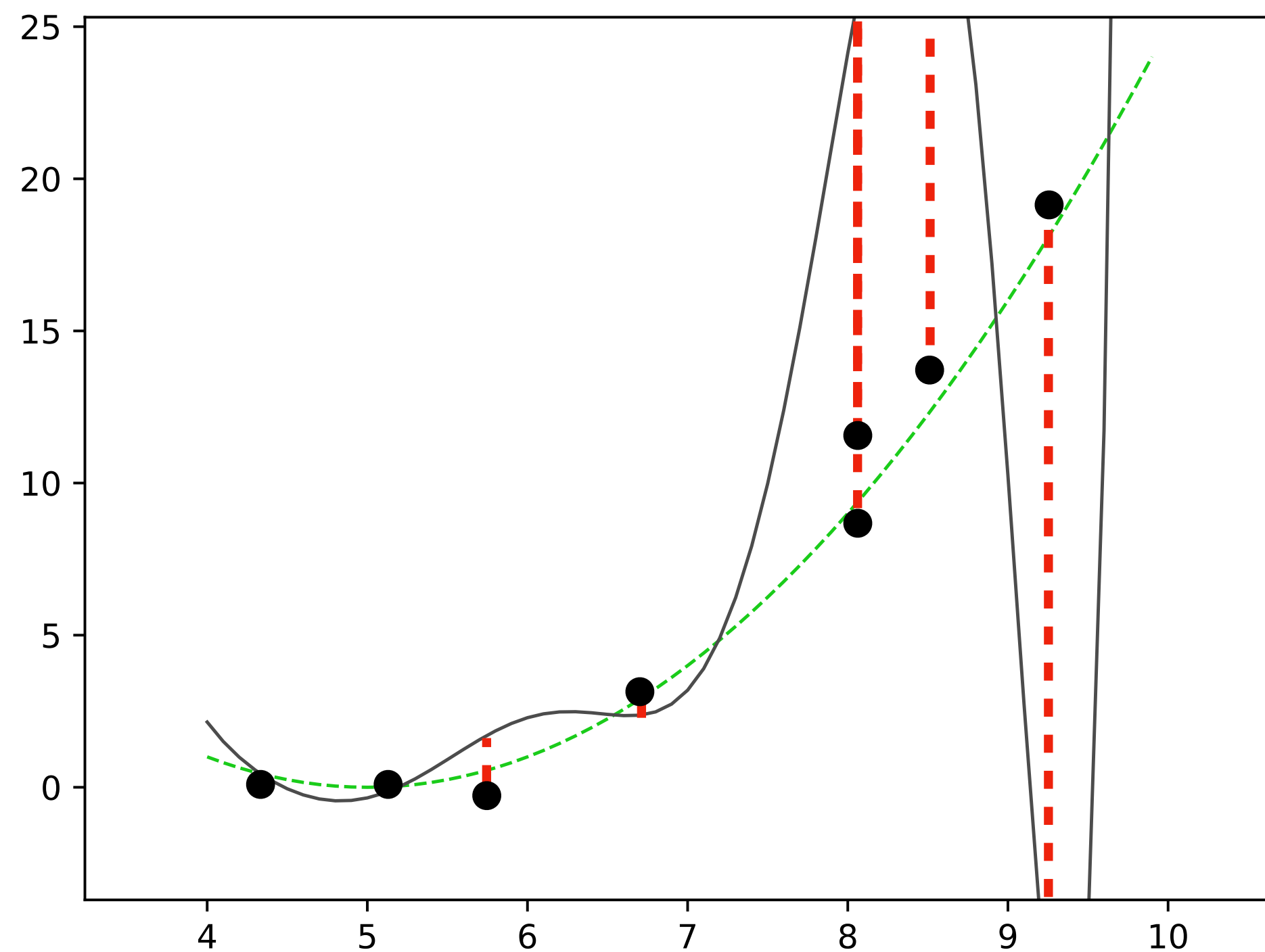
Training error versus generalization error



["Deep Learning", Goodfellow et al.]

How do we know if we are underfitting or overfitting?

Validation data $\{x_i^{(\text{val})}, y_i^{(\text{val})}\}$



Cross validation: measure prediction error on validation set

Fitting a model

Underfitting?

1. add more parameters (more features, more layers, etc.)

Overfitting?

1. remove parameters
2. add **regularizers**

Regularization

Empirical risk minimization:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + R(\theta)$$

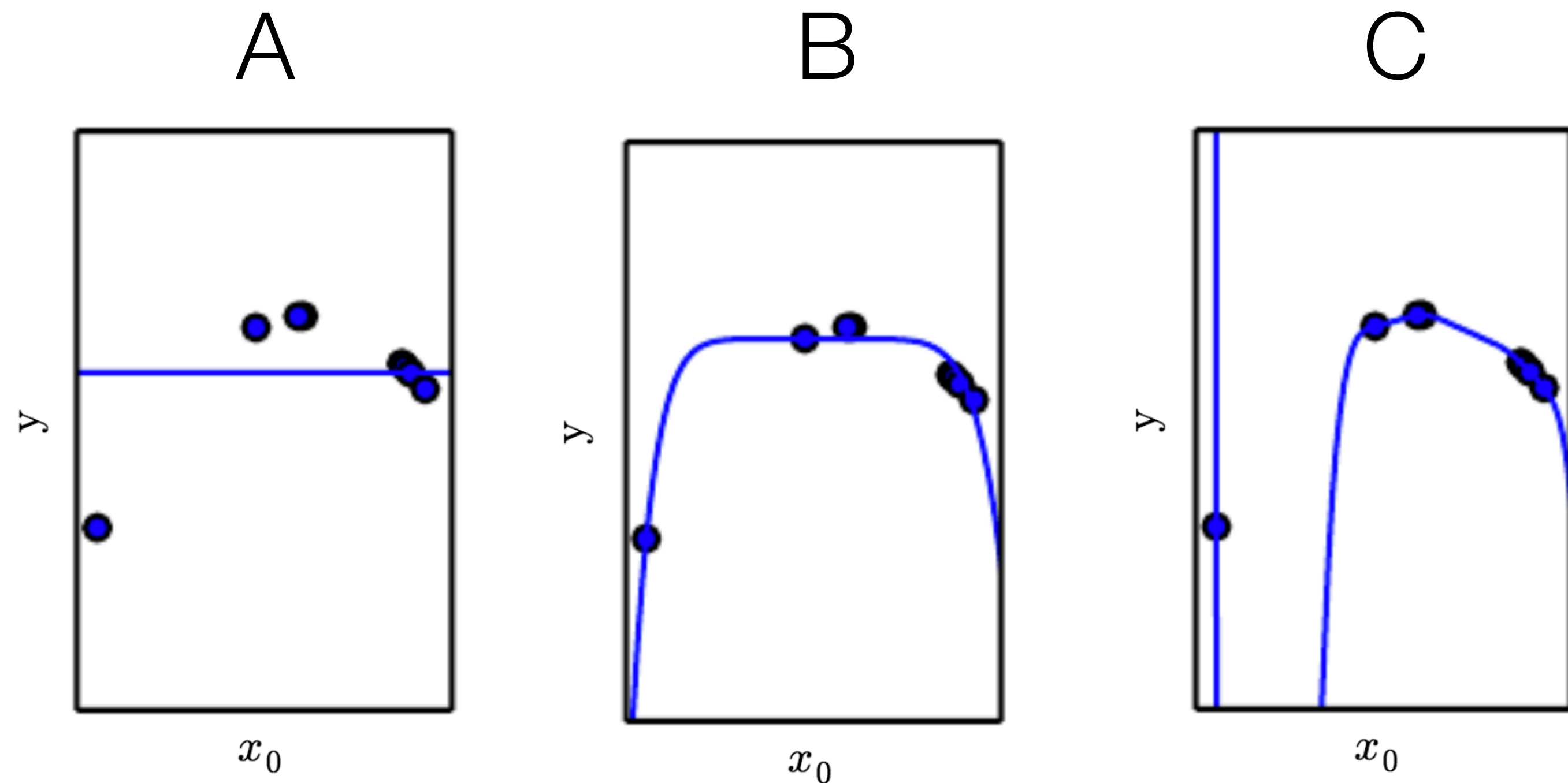
Regularized least squares

$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

$$R(\theta) = \lambda \|\theta\|_2^2 \longleftarrow \text{Only use polynomial terms if you really need them! Most terms should be zero}$$

ridge regression, a.k.a., **Tikhonov regularization**

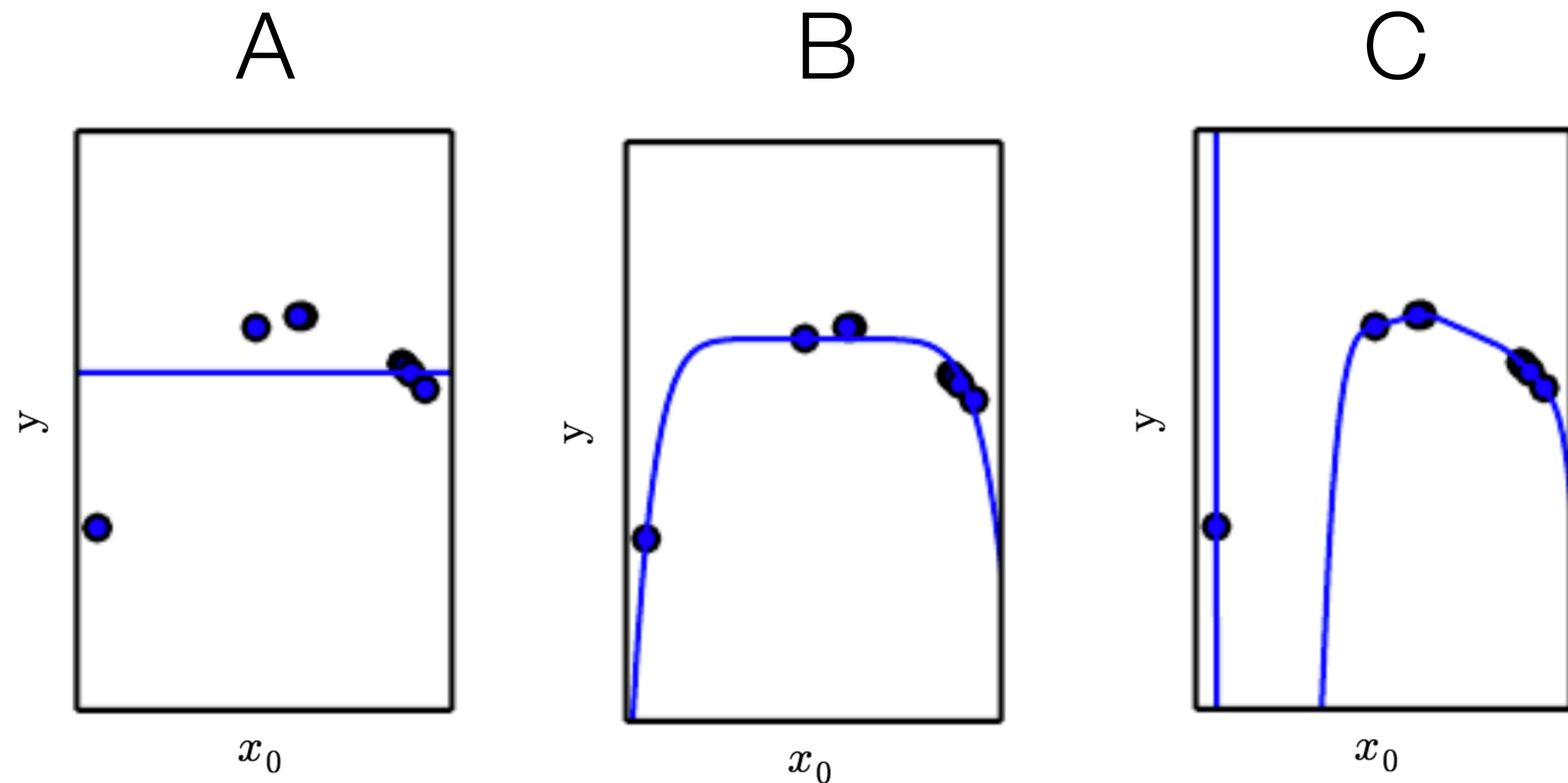
$$\theta^* = \arg \min_{\theta} \sum_{I=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + \lambda \|\theta\|_2^2$$



Low λ – ?
Medium λ – ?
High λ – ?

[Adapted from “Deep Learning”, Goodfellow et al.]

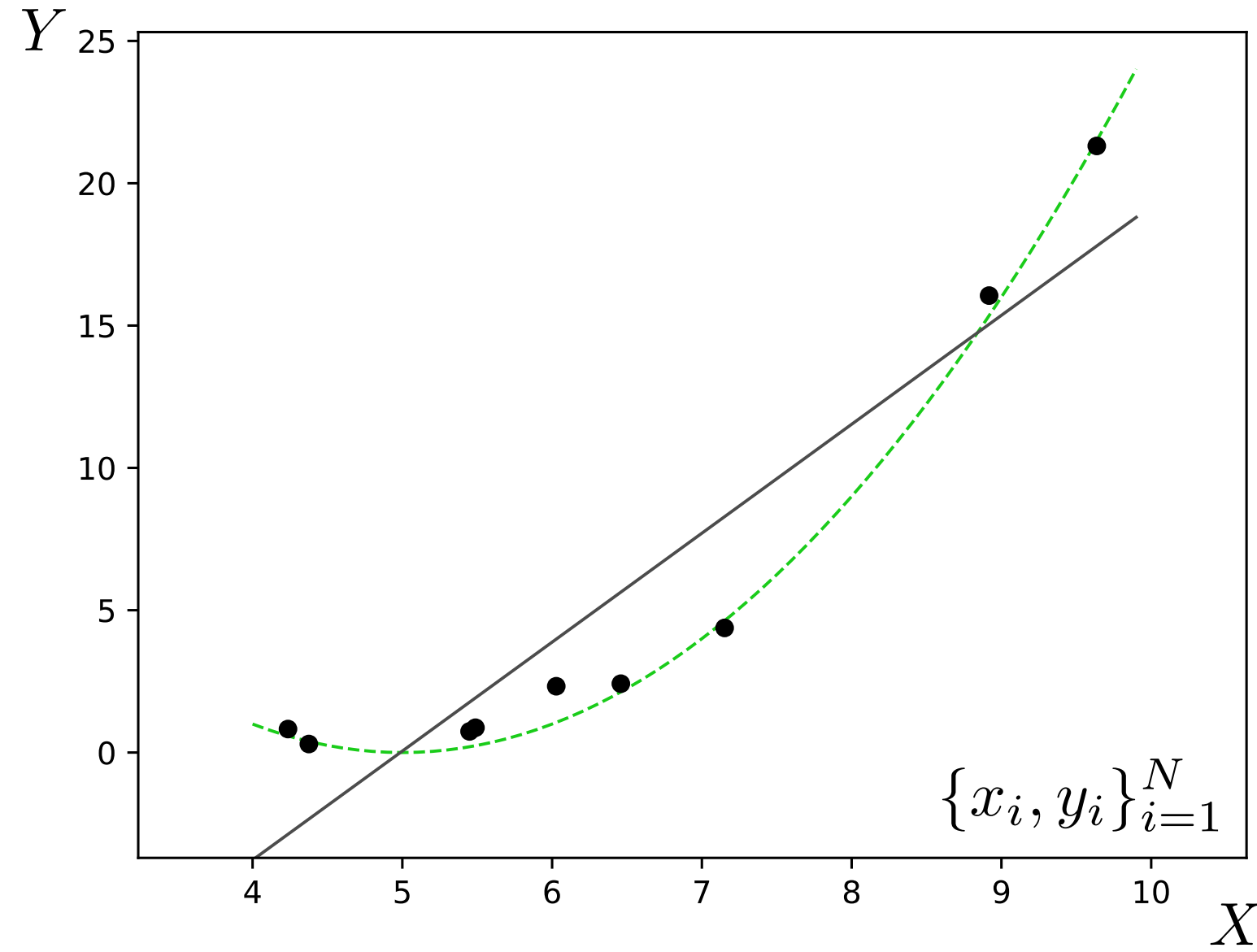
$$\theta^* = \arg \min_{\theta} \sum_{I=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + \lambda \|\theta\|_2^2$$



Low λ – C
Medium λ – B
High λ – A

Underfitting

$$K = 1$$

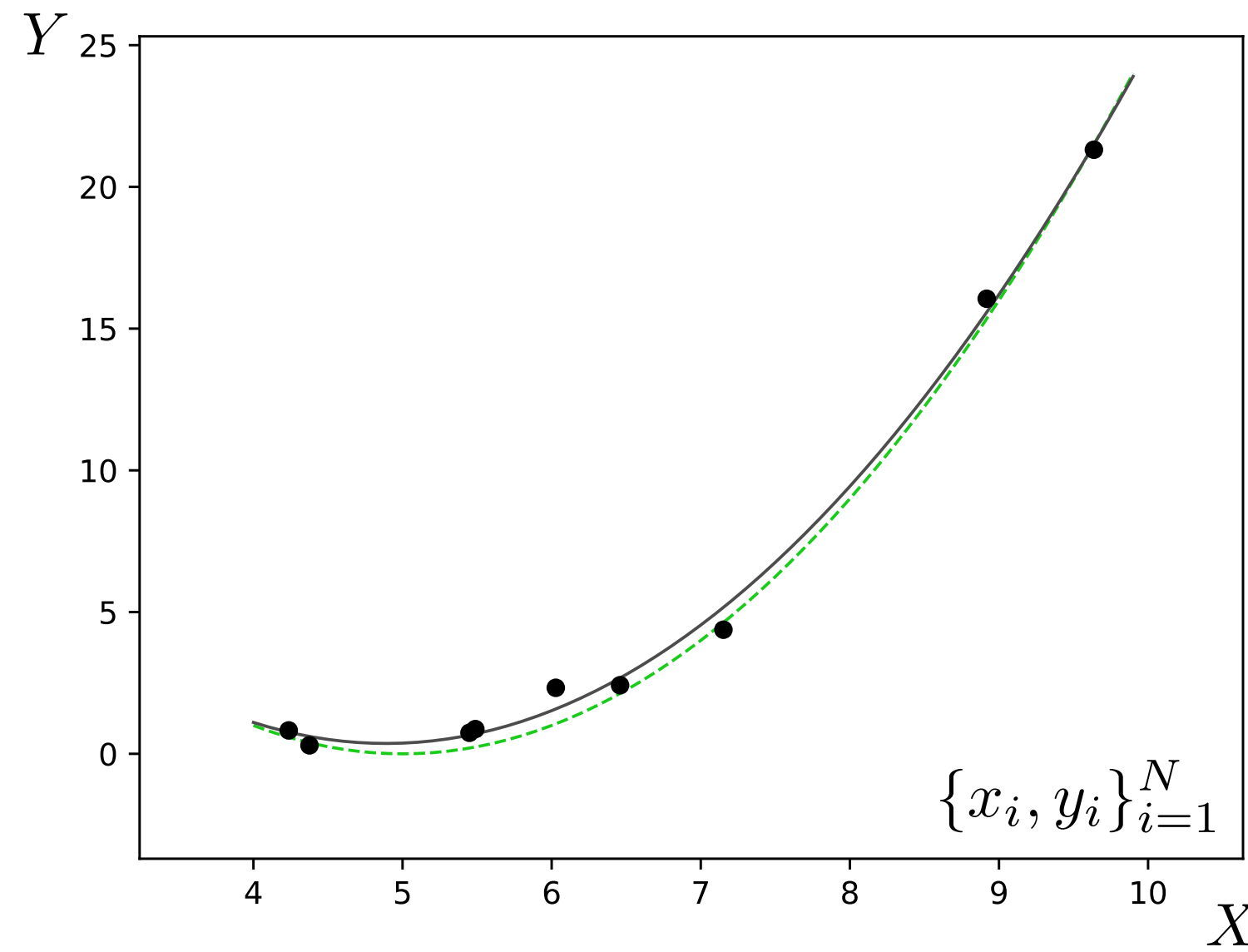


Simple model

Doesn't fit the training data

Appropriate model

$$K = 2$$

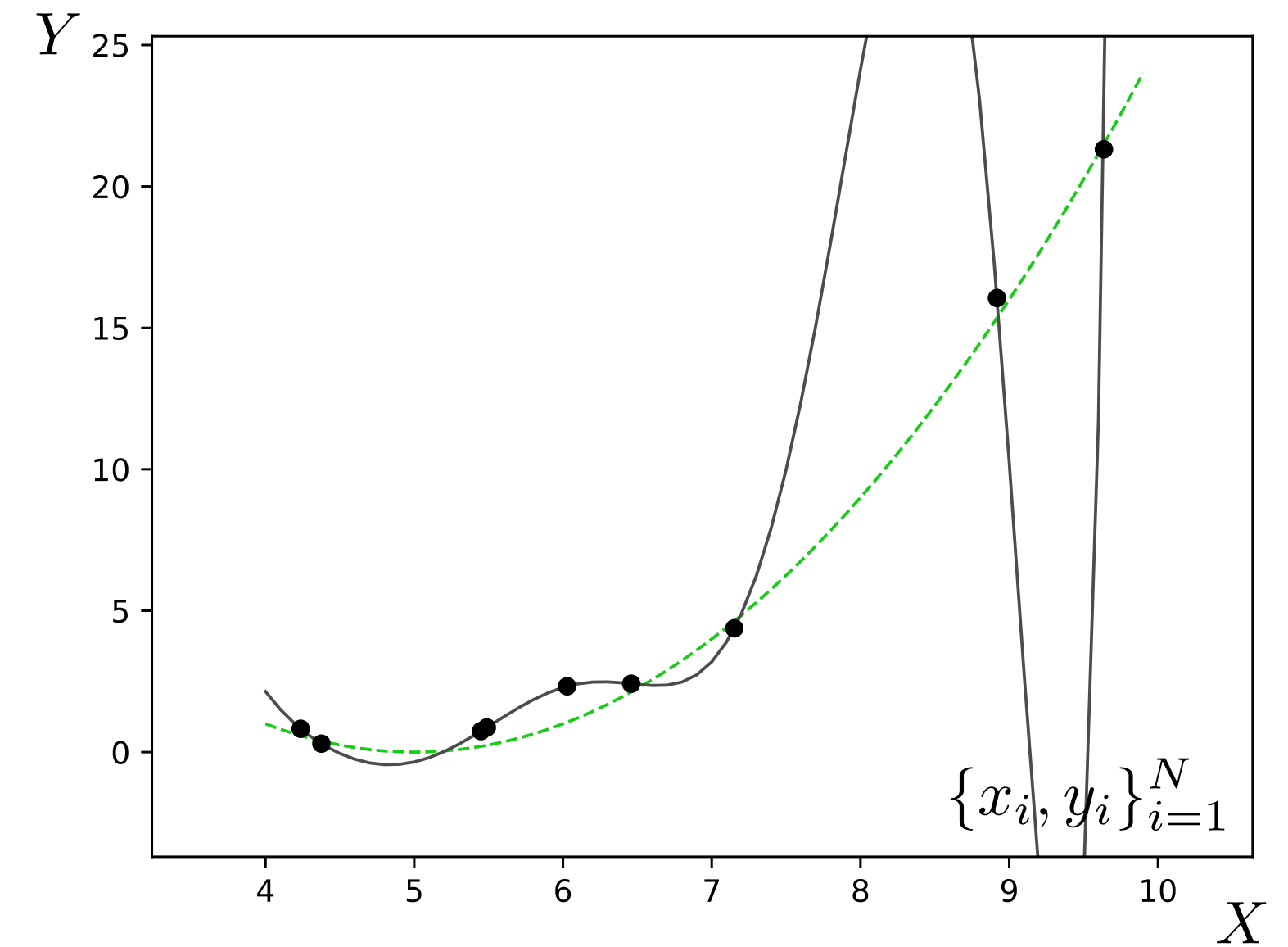


Simple model

Fits the training data

Overfitting

$$K = 10$$



Complex model

Fits the training data

Image classification

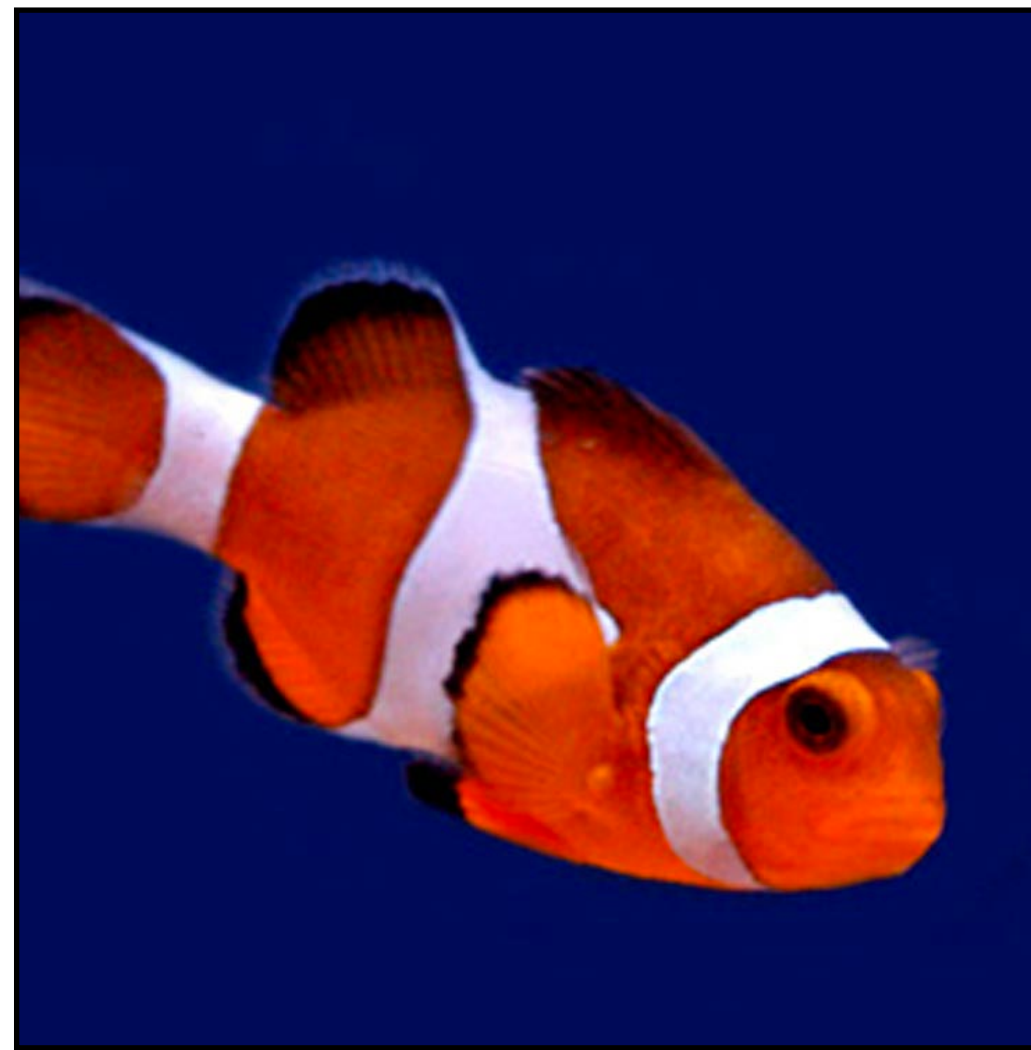


image x



"Fish"

label y

Image classification



image \mathbf{x}



"Fish"

label \mathbf{y}

Image classification



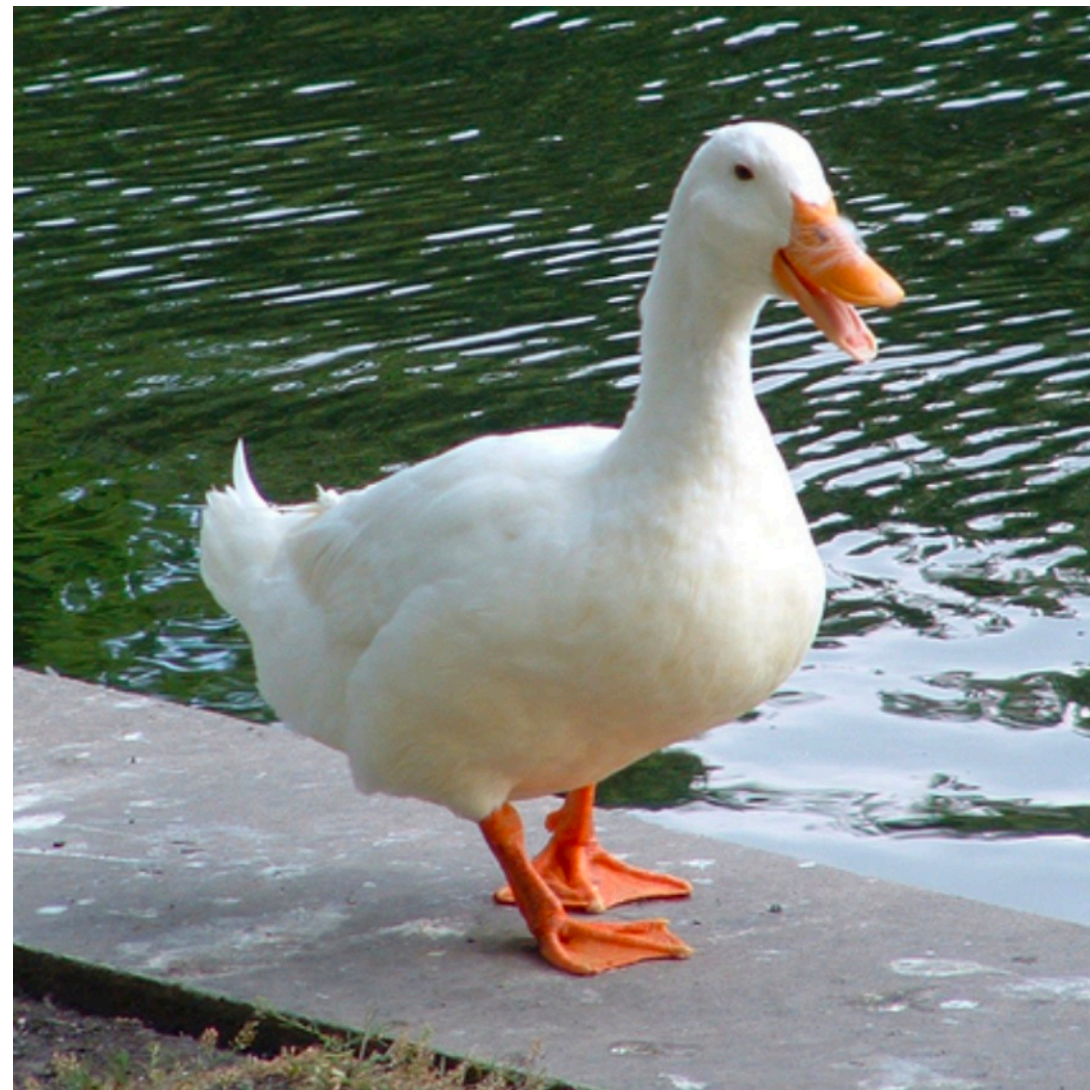
image x



"Fish"

label y

Image classification



⋮

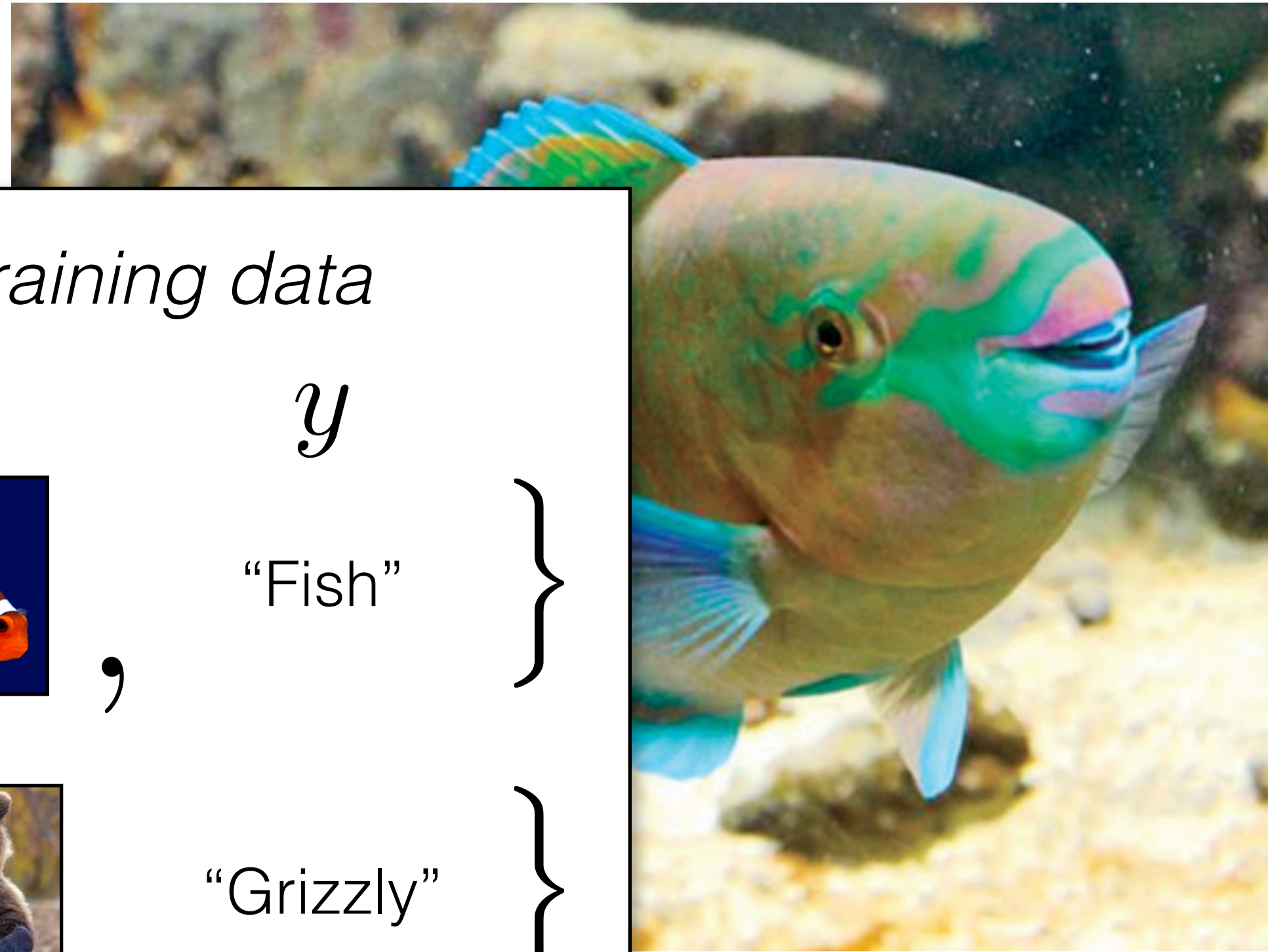
image \mathbf{x}



“Duck”

label y

\mathbf{x}



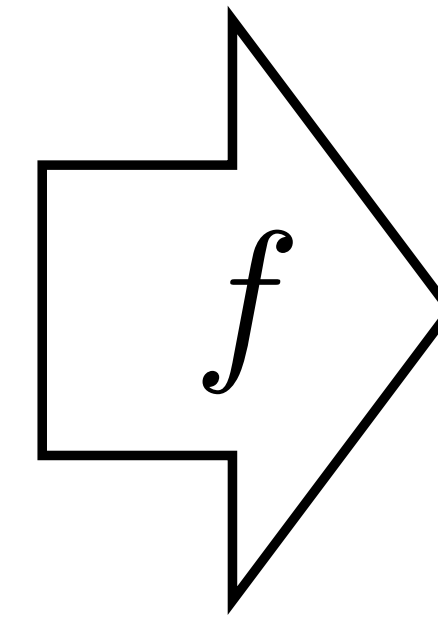
Training data

\mathbf{x}

y



⋮



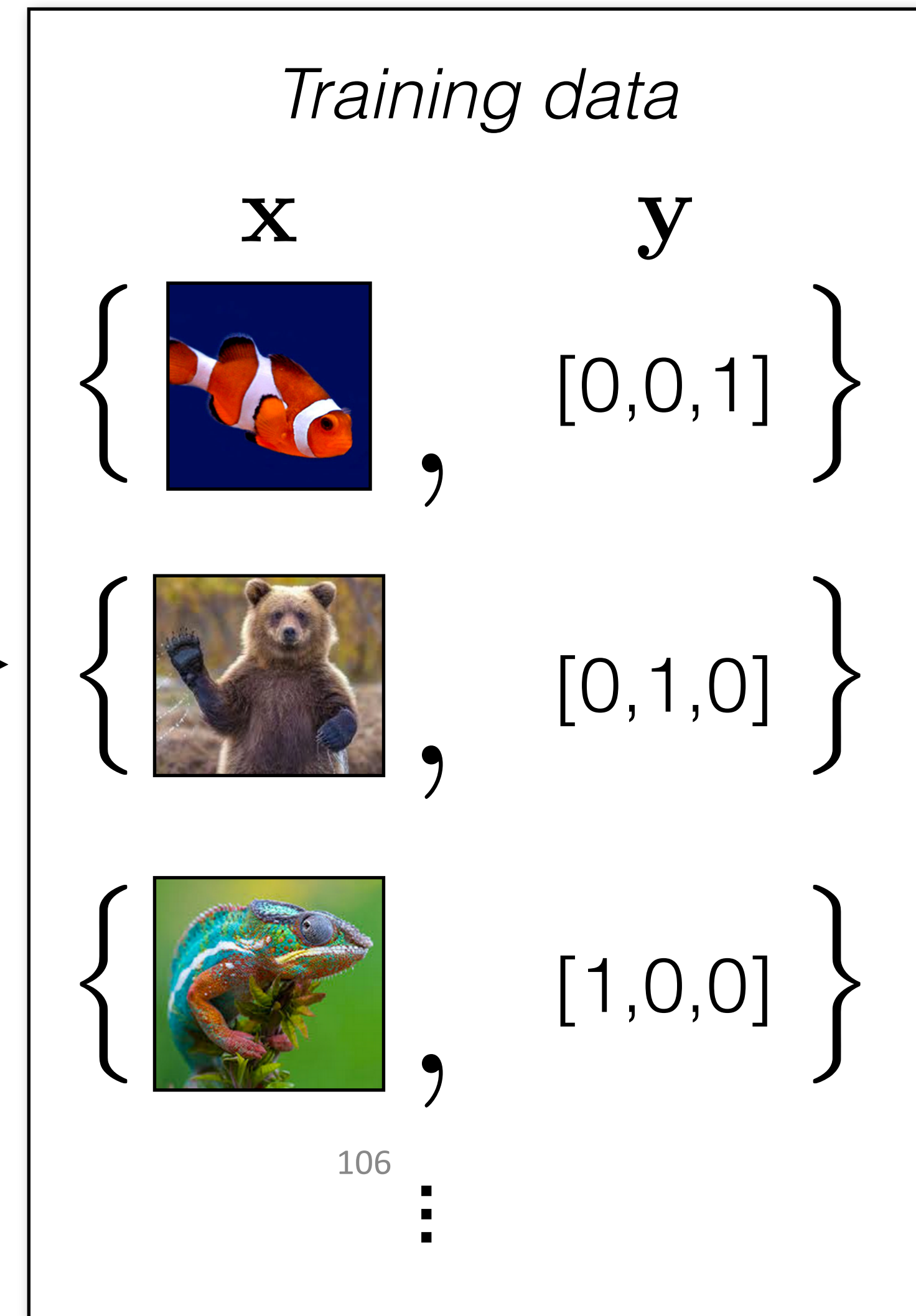
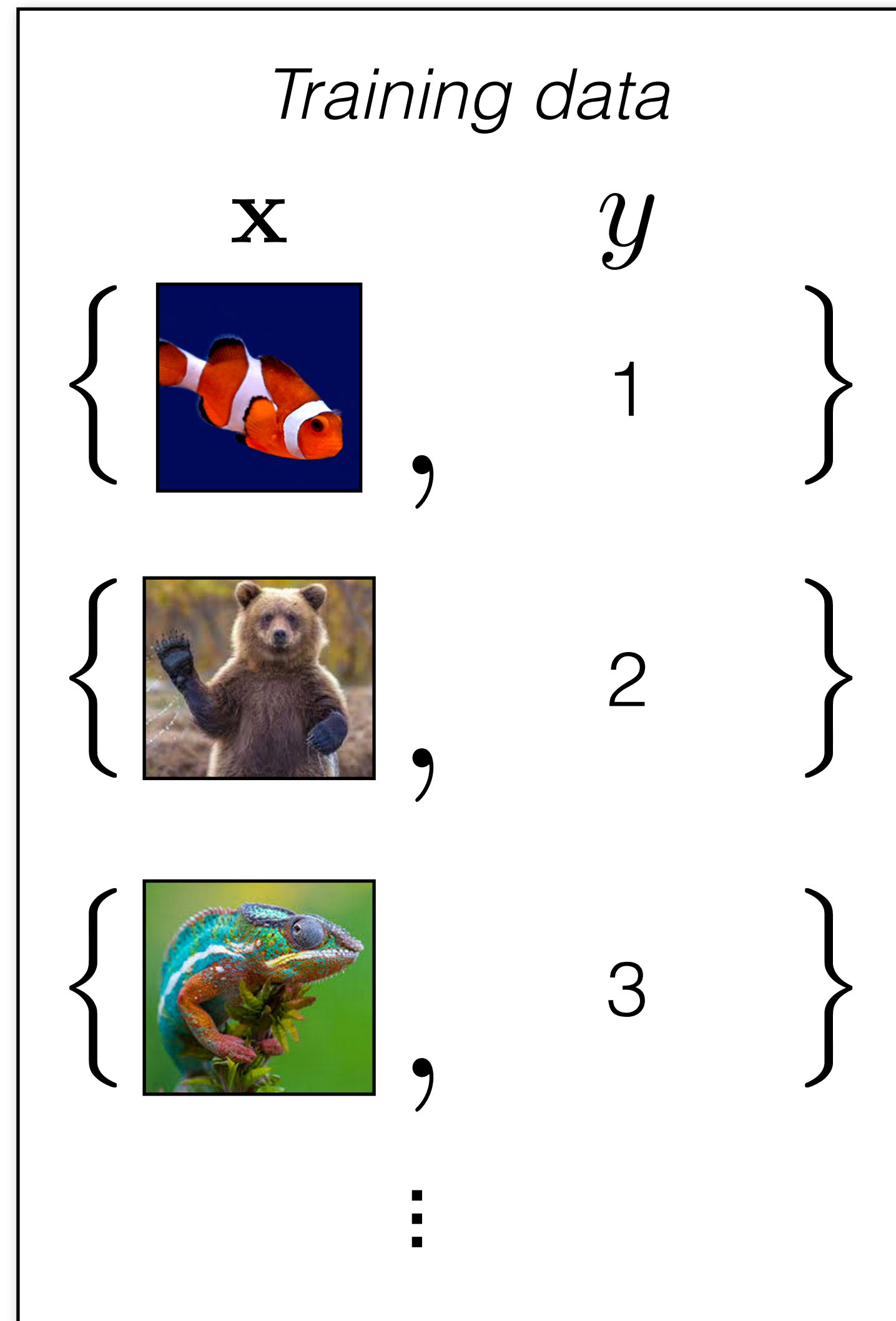
y

“Fish”

$$\arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i), y_i)$$

How to represent class labels?

One-hot vector



Recall: loss function

0-1 loss (number of misclassifications)

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = 1 - \mathbb{1}(\hat{\mathbf{y}}, \mathbf{y}) \leftarrow \text{number of misclassifications}$$

Least squares (predict 1 for true class, 0 for others)

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{k=1}^K (y_k - \hat{y}_k)^2$$

Cross entropy: a good surrogate that we'll be able to optimize:

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

Ground truth label y

x

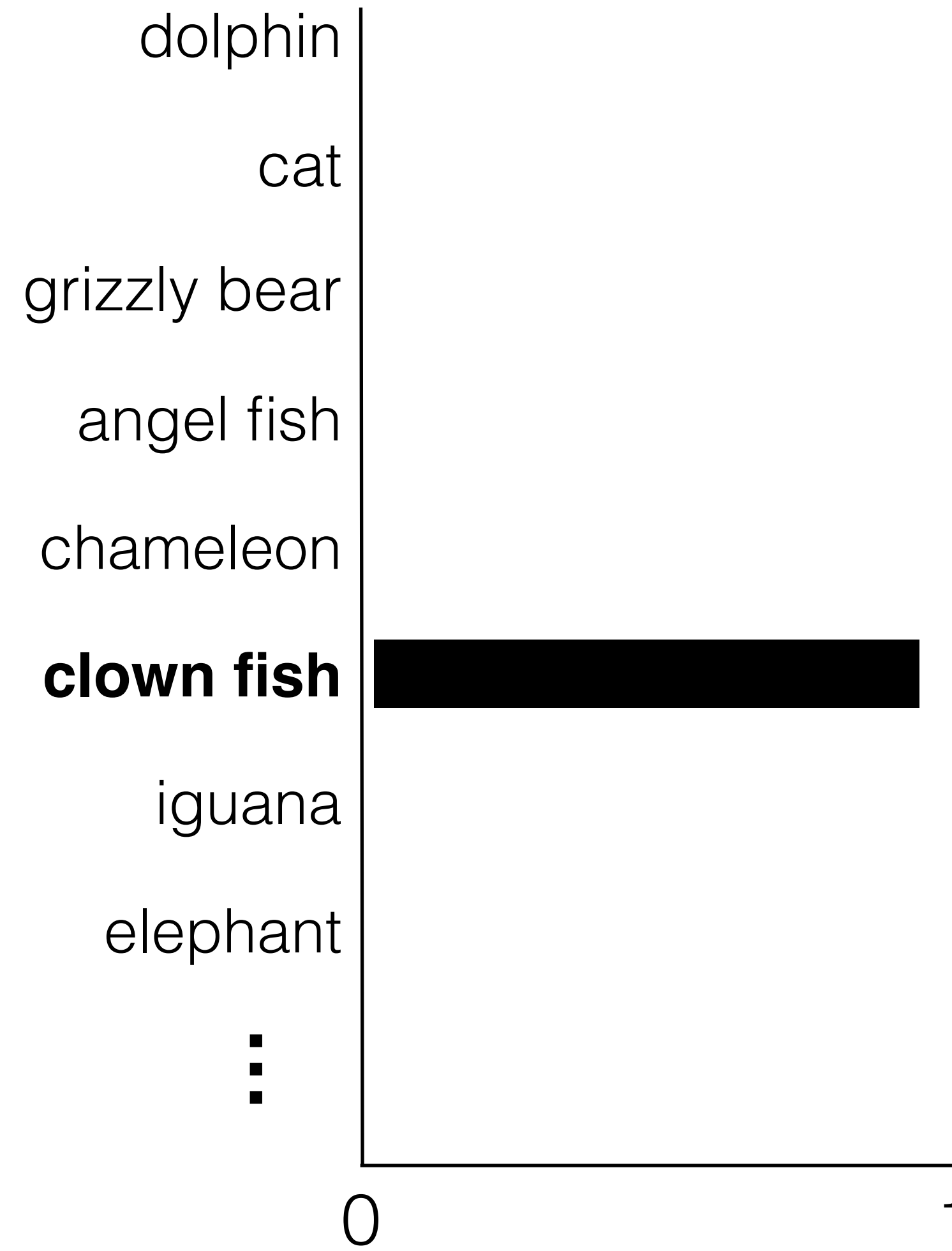


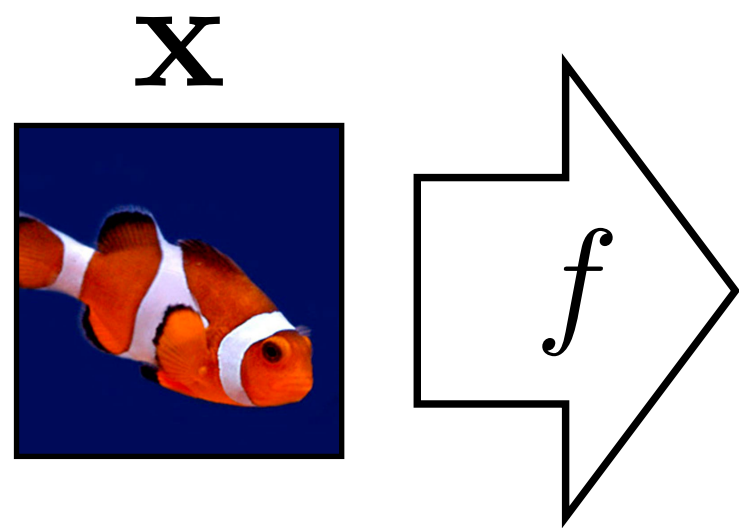
$[0, 0, 0, 0, 0, 1, 0, 0, \dots]$



X

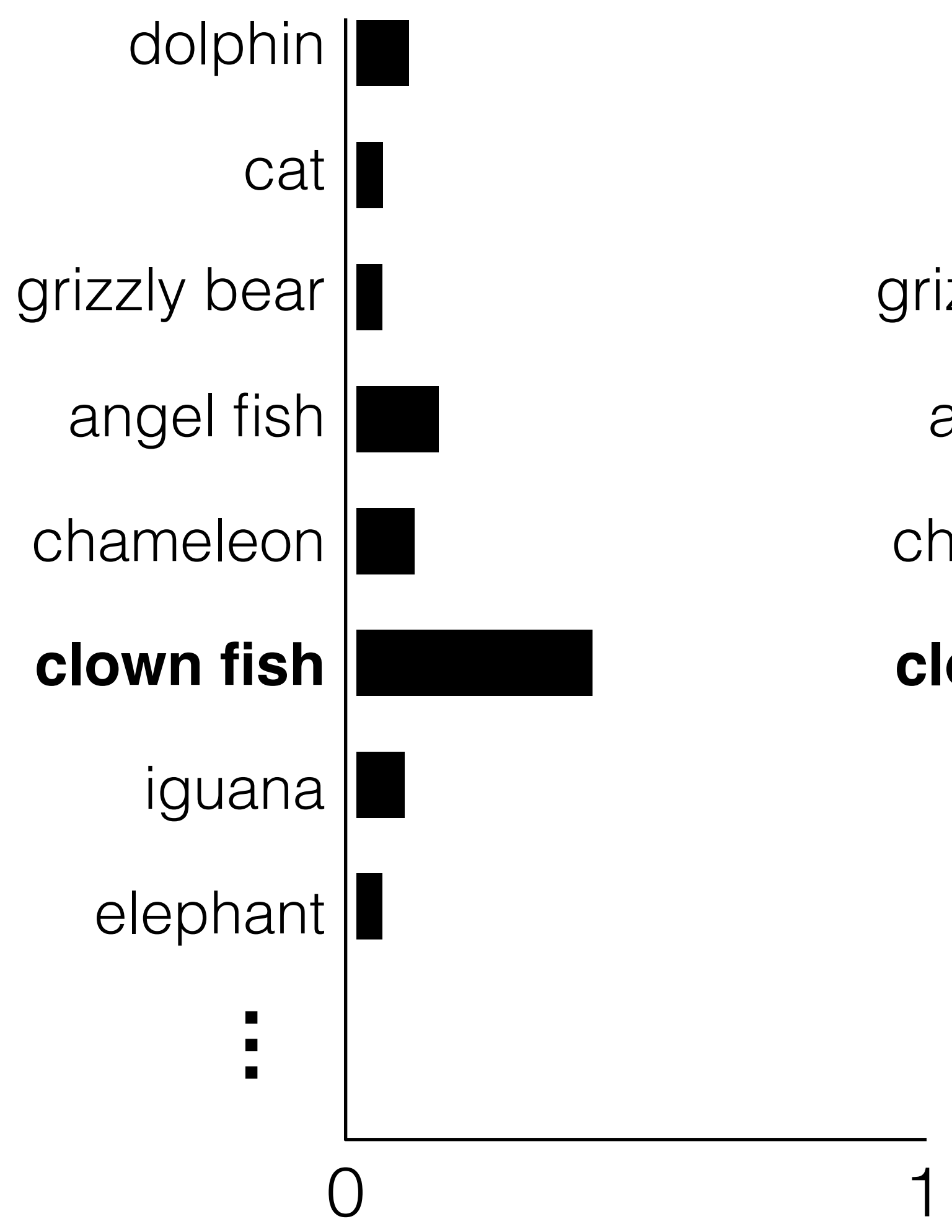
Ground truth label **y**



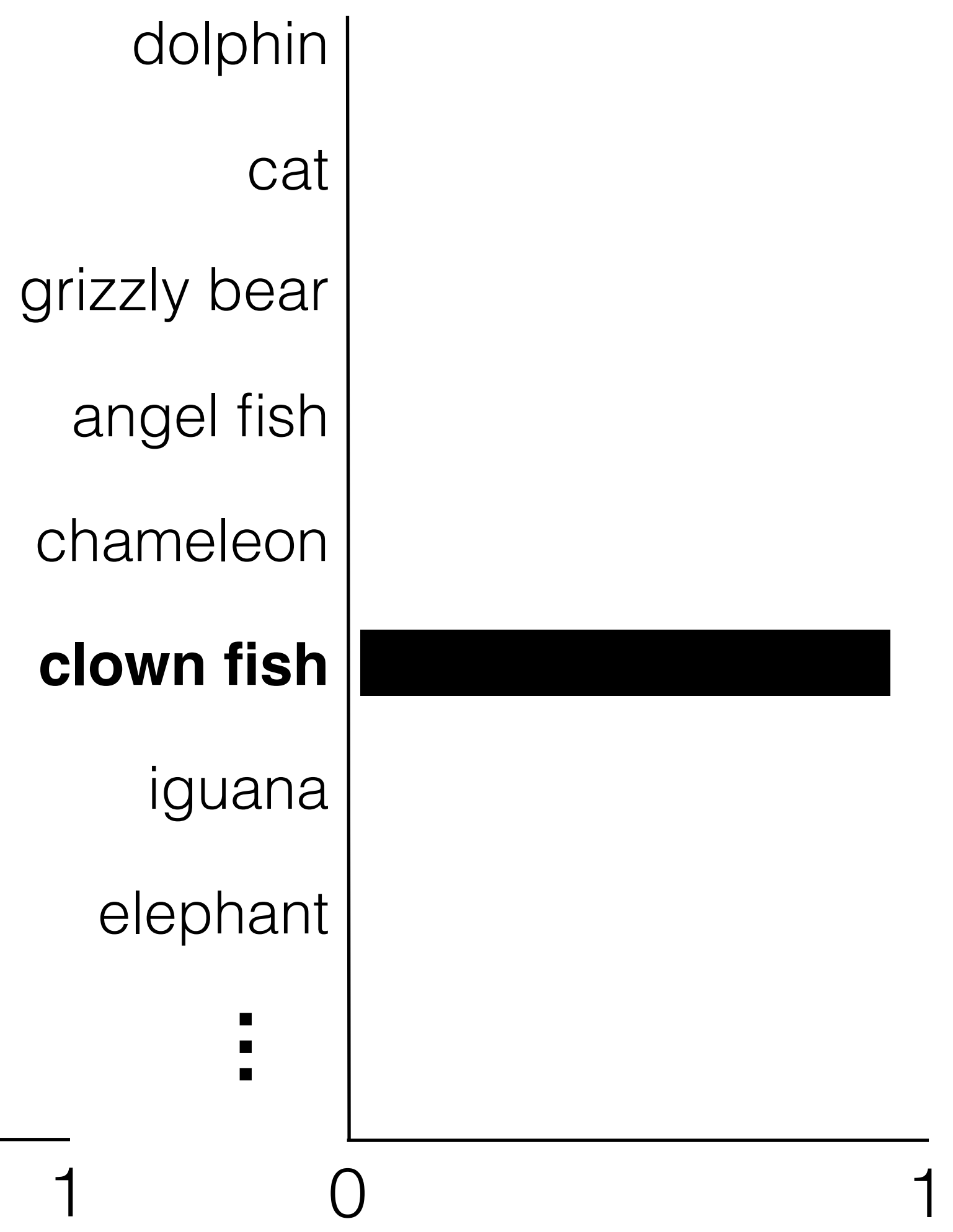


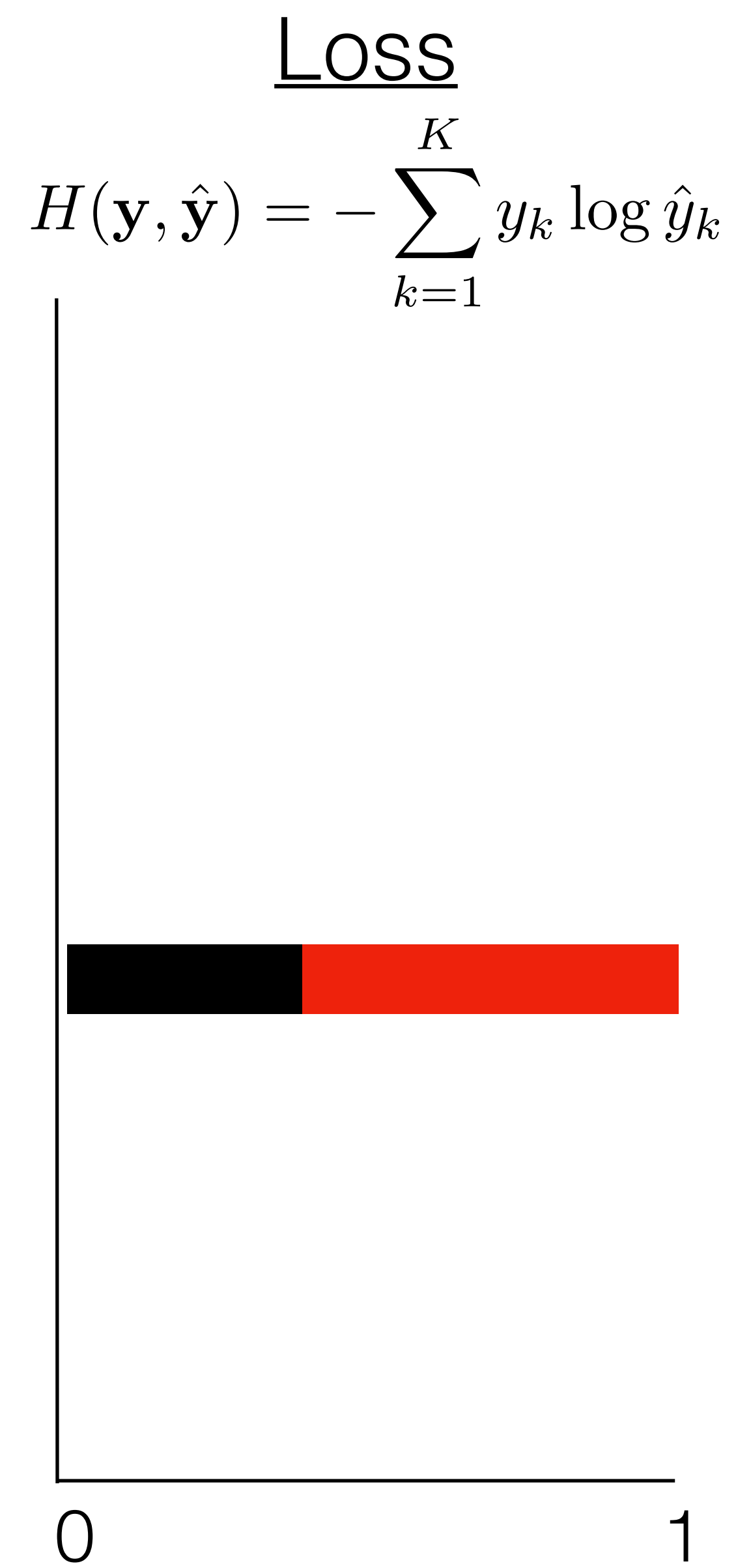
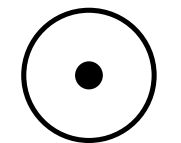
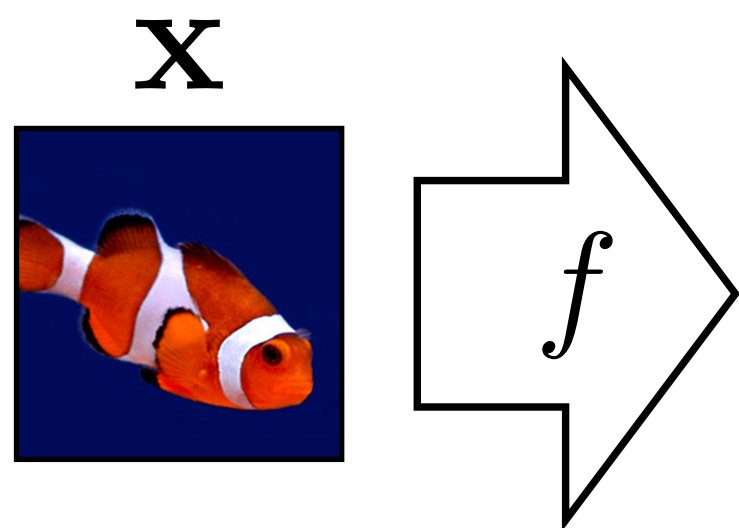
Prediction \hat{y}

$$f_{\theta} : X \rightarrow \mathbb{R}^K$$



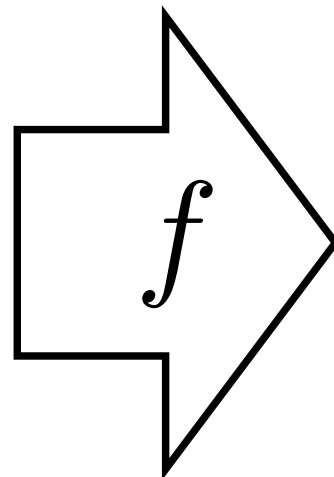
Ground truth label y





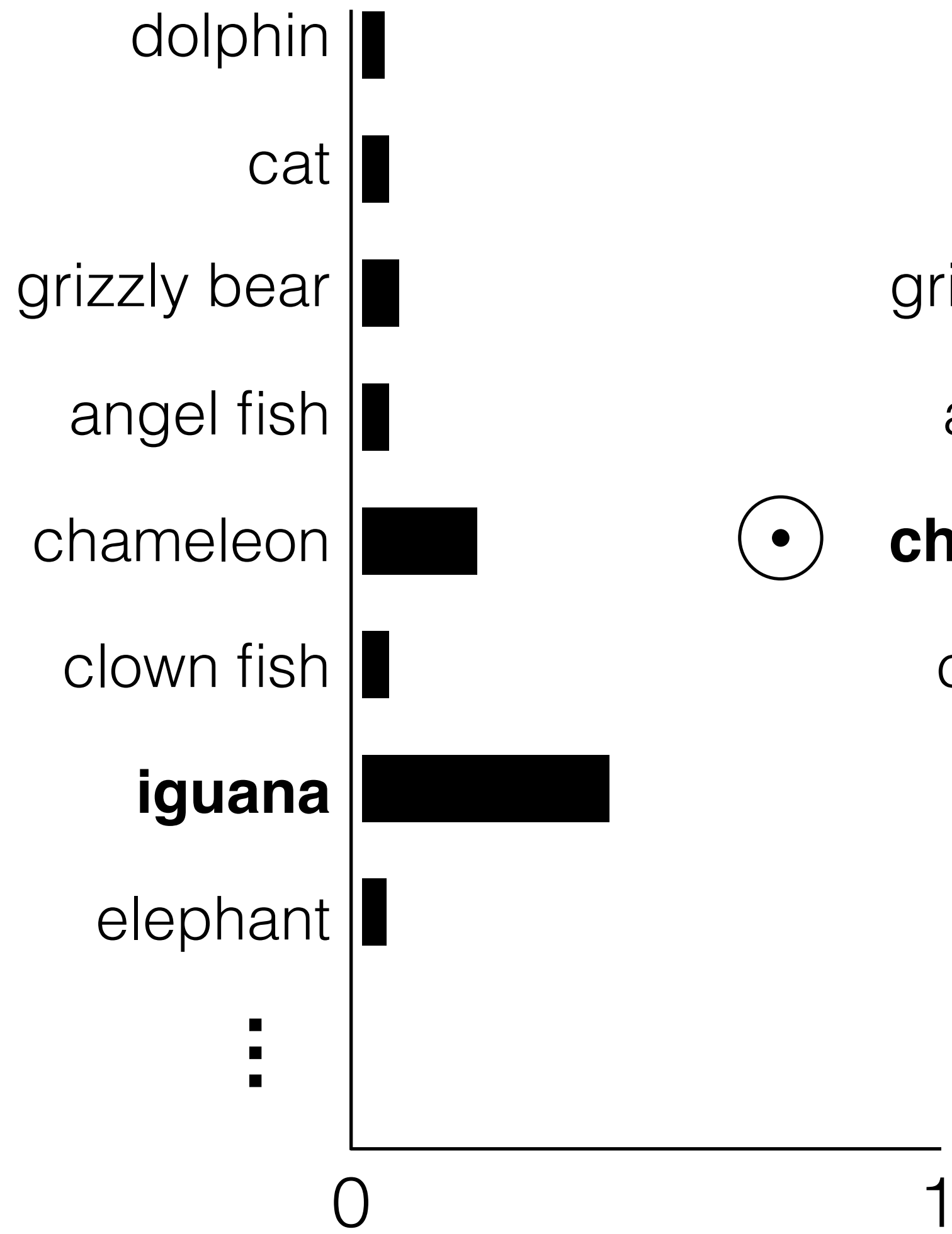


X

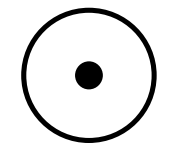


Prediction \hat{y}

$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

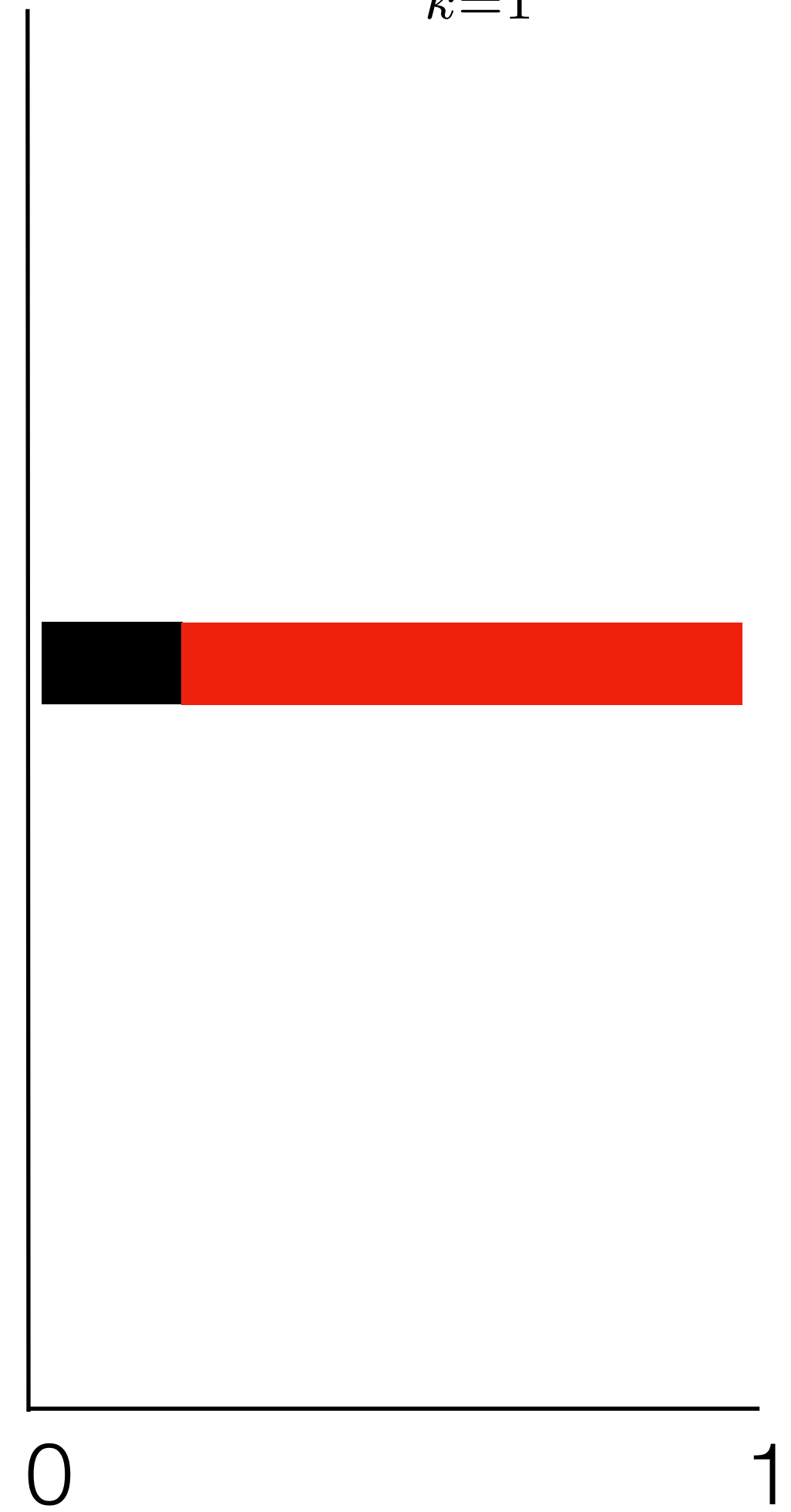


Ground truth label y



Loss

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$



Softmax regression (a.k.a. multinomial logistic regression)

$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

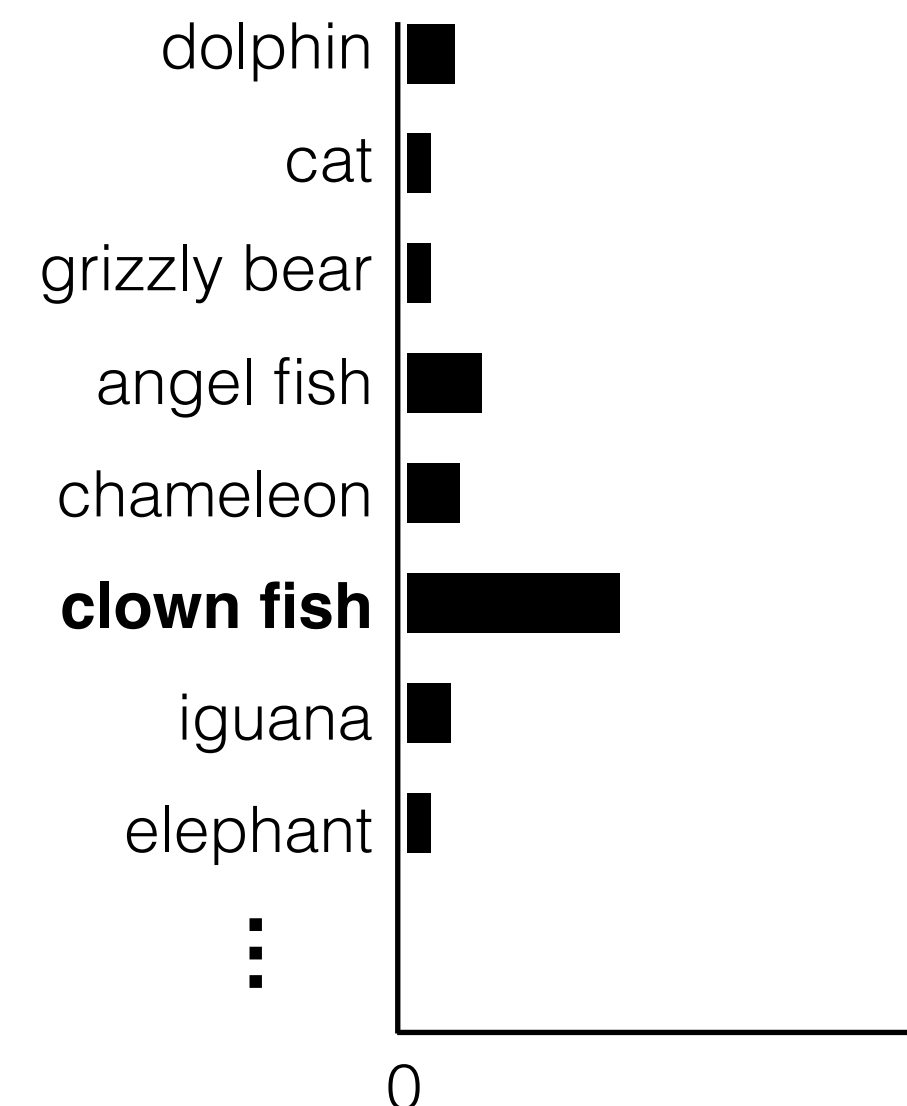
← **logits**: vector of K scores, one for each class

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

← squash into a non-negative vector that sums to 1
— i.e. **a probability density function**

$$\hat{y}_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

$\hat{\mathbf{y}} =$



Next lecture: more linear models