

# Lecture 6: Machine learning

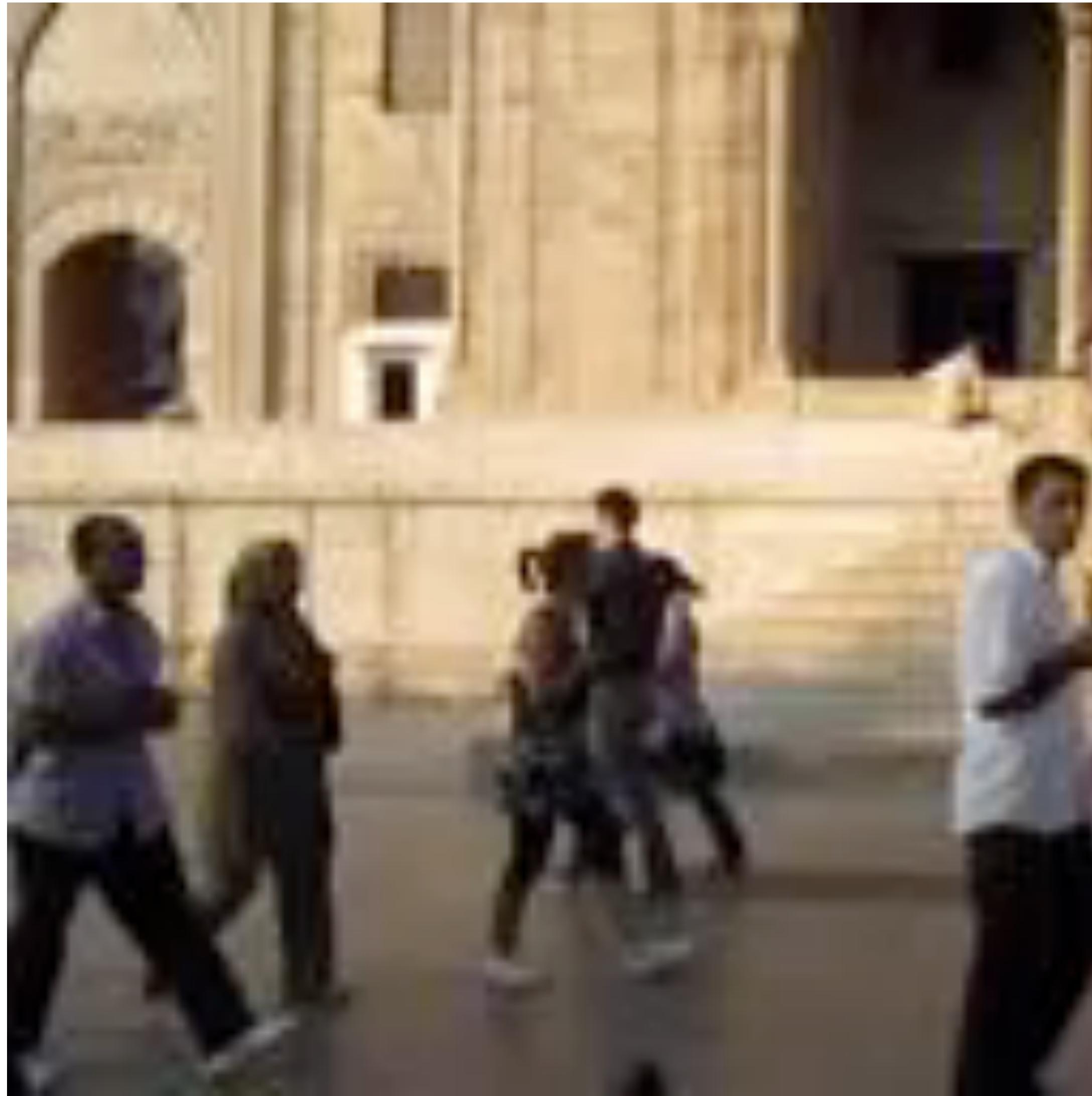
- My office hours cancelled this week
  - Can do them by appointment instead
- PS3 out tonight

# Today

- Temporal filtering
- Introduction to machine learning

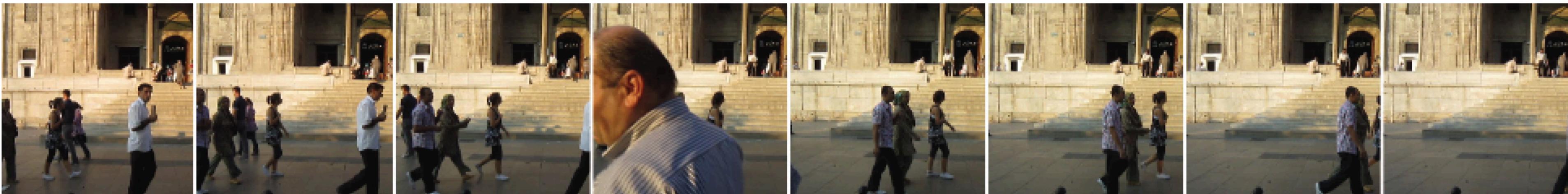
# Temporal filtering

# Temporal filtering

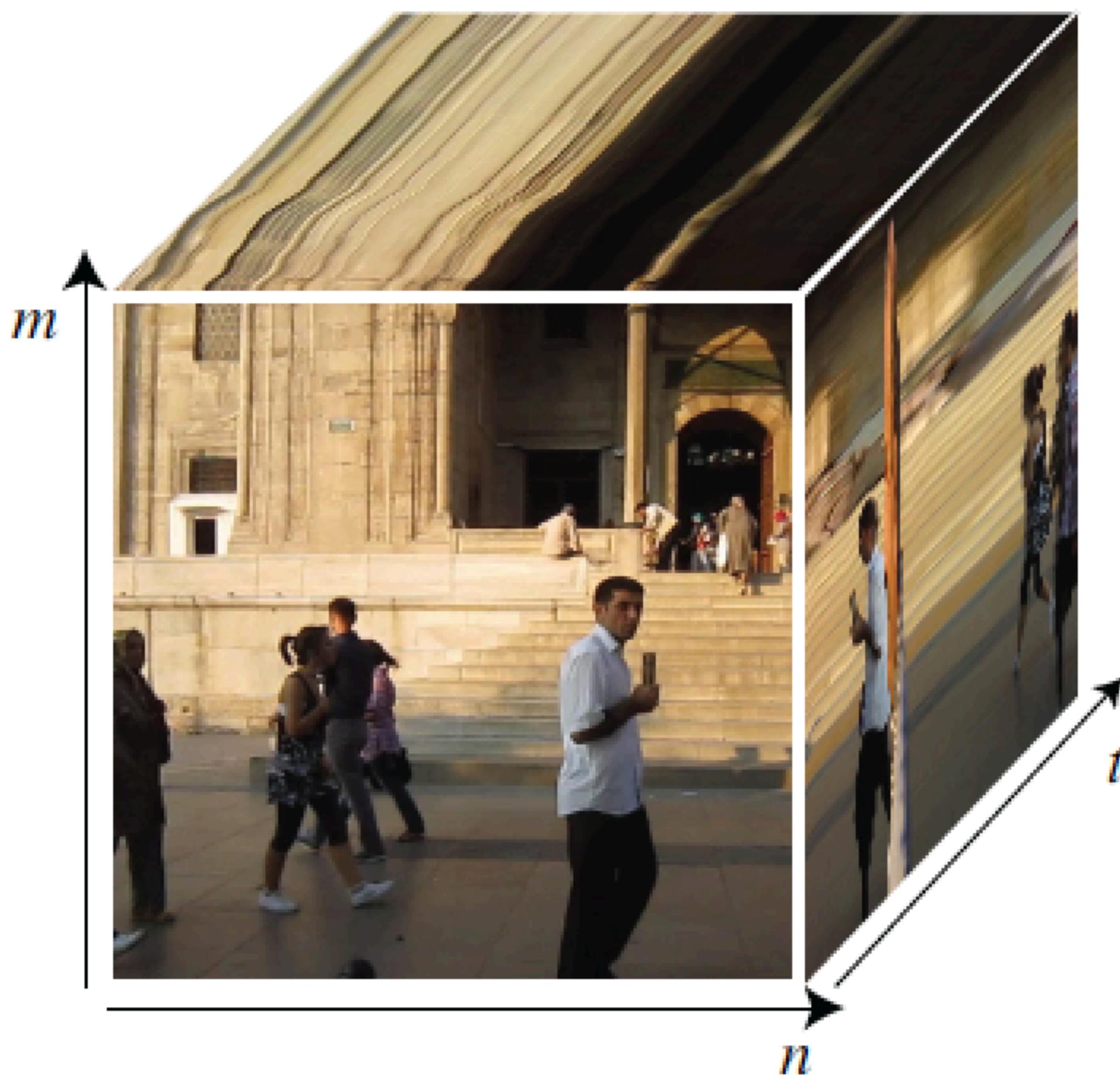


Source: Freeman, Torralba, Isola

# Videos

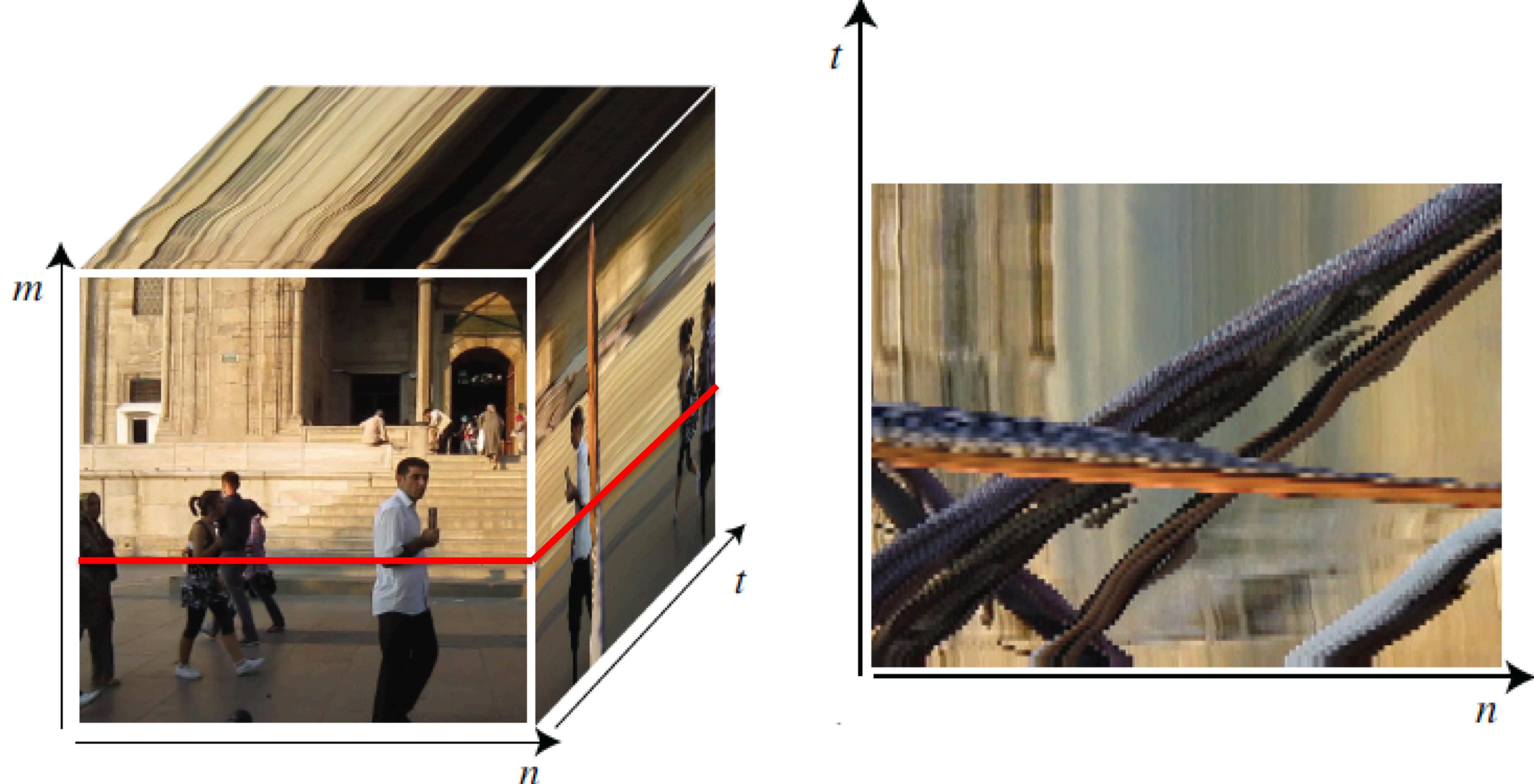


time →



Source: Freeman, Torralba, Isola

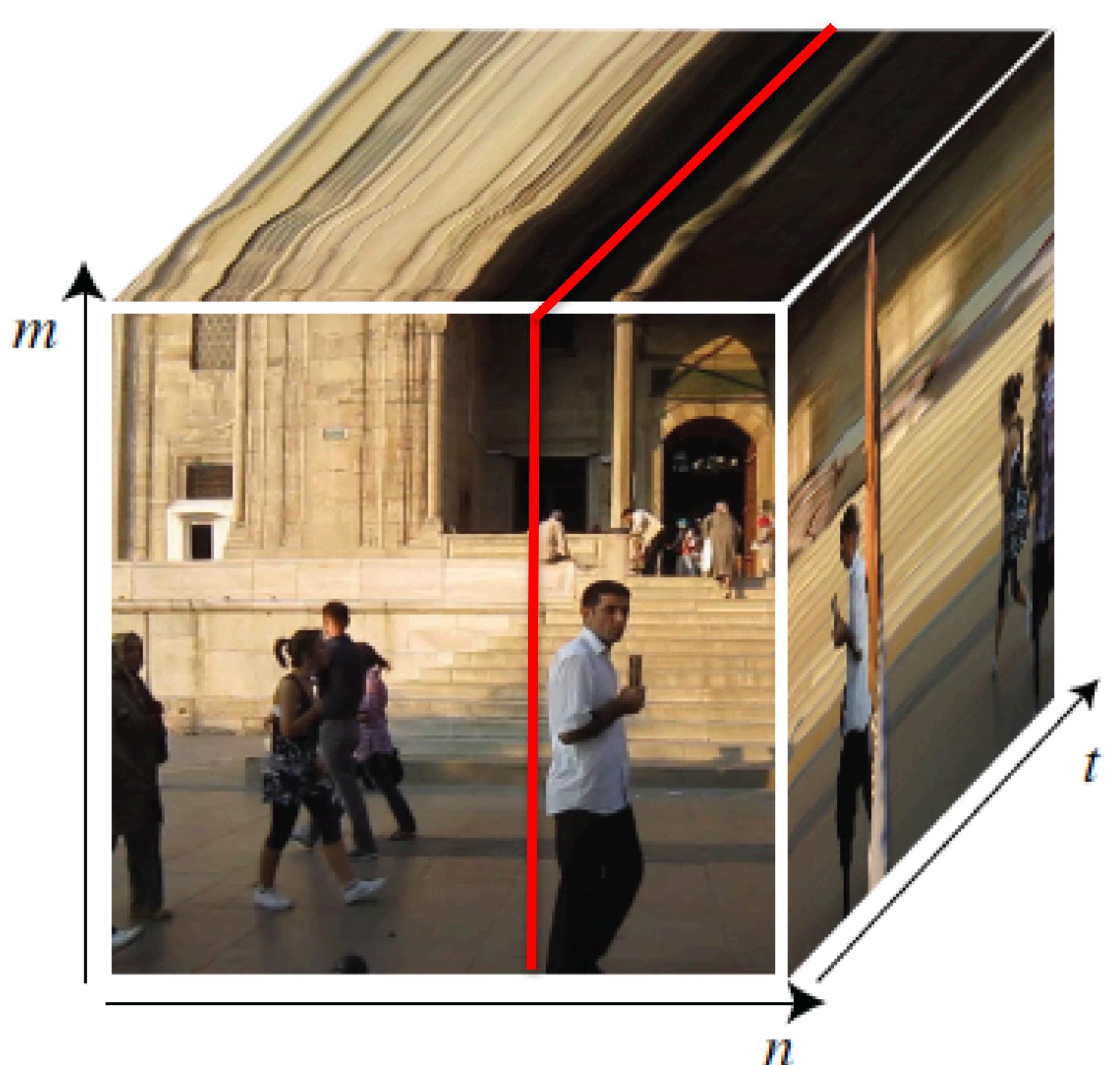
# Videos



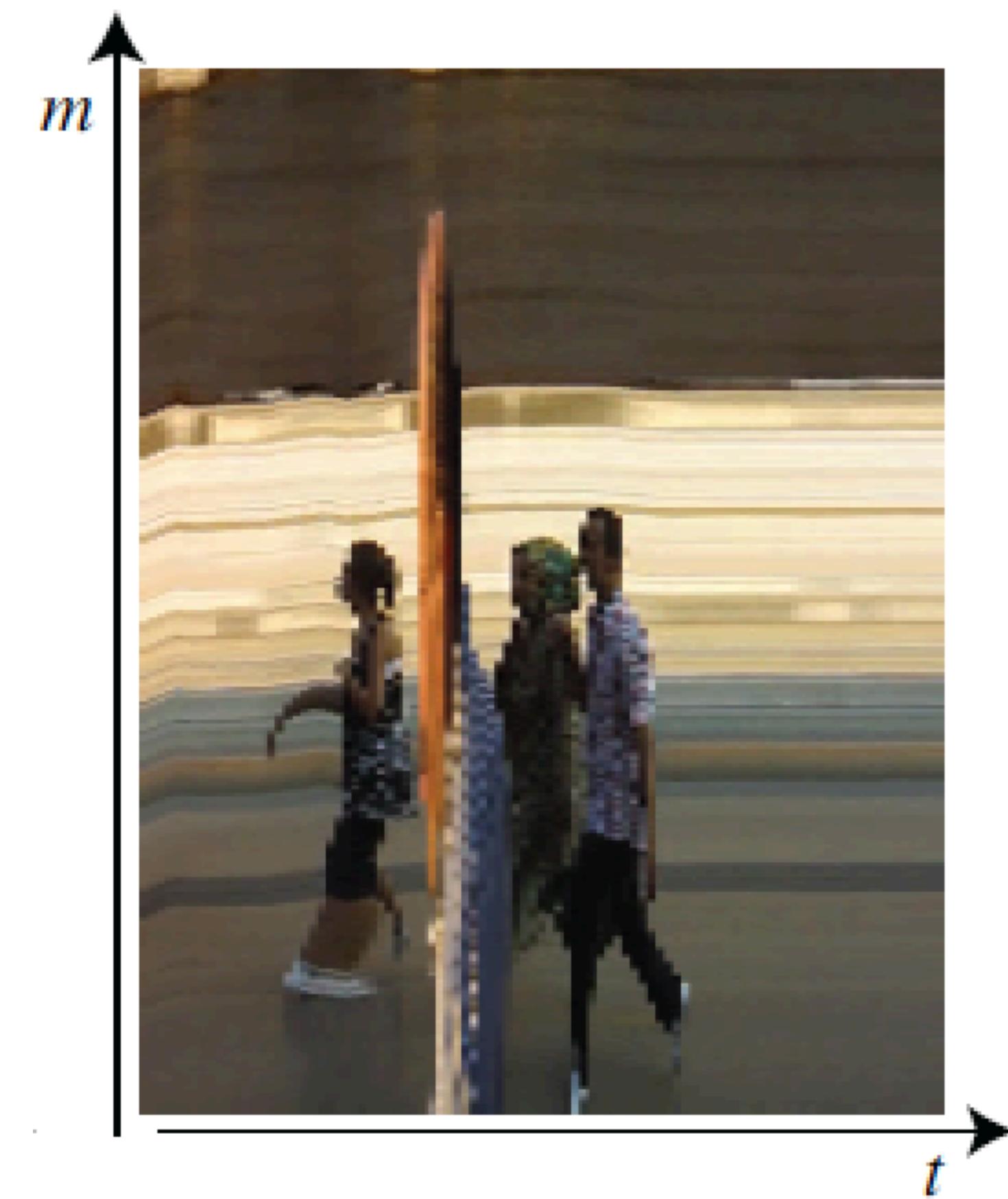
Cube size = 128x128x90

Source: Freeman, Torralba, Isola

# Videos

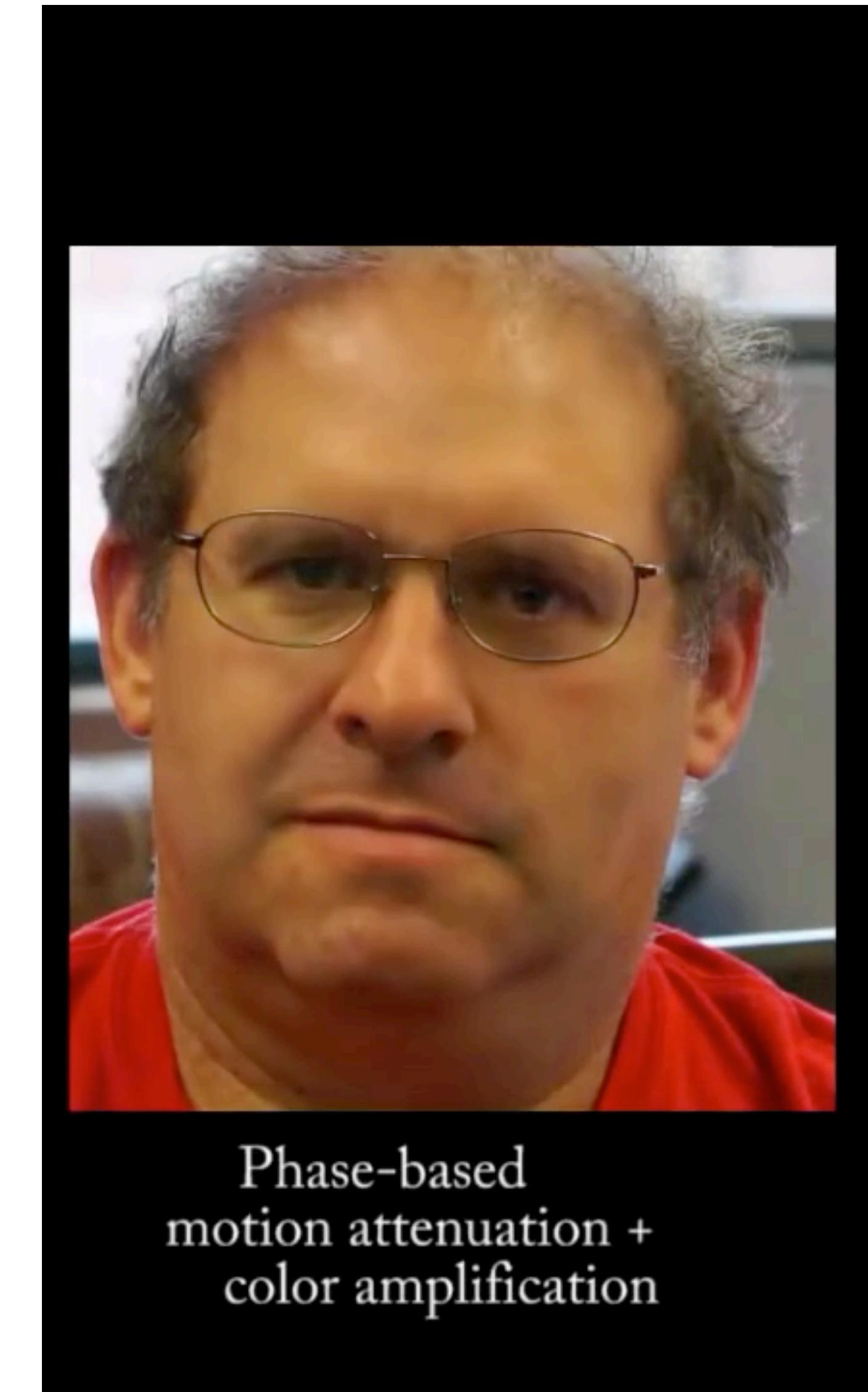


Cube size = 128x128x90



Source: Freeman, Torralba, Isola

# Color amplification



# Motion magnification

Original



- Just spatial-temporal filtering!
- You'll do a simplified version in PS3

Original

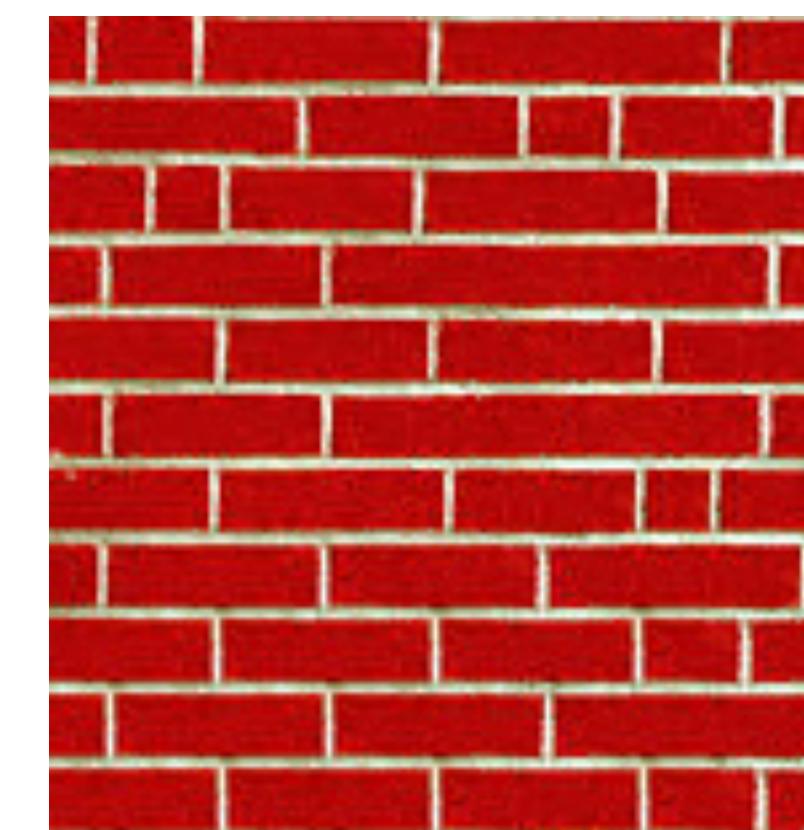
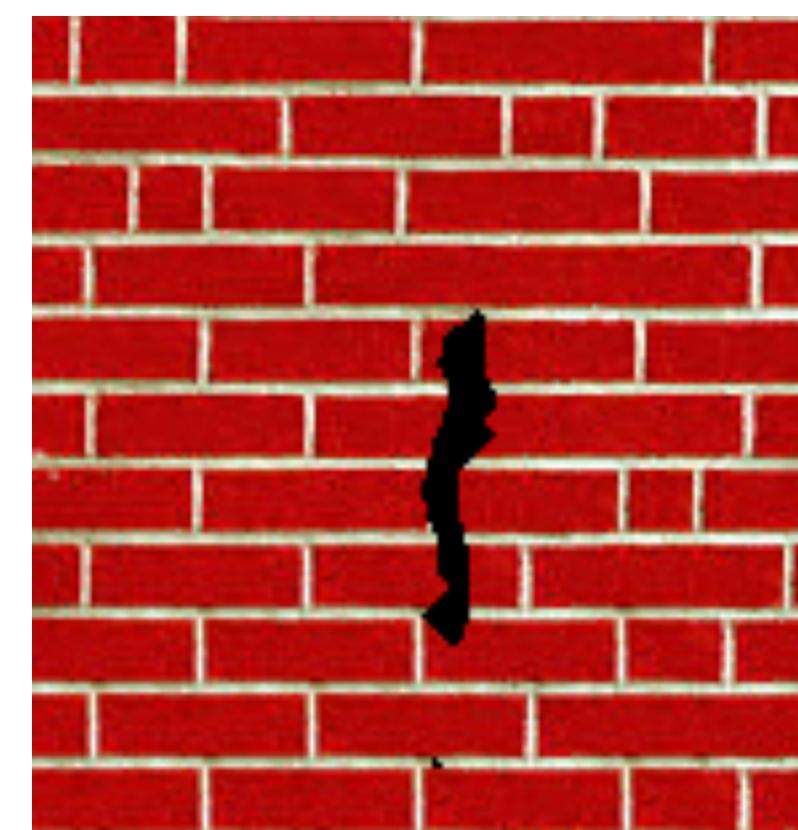


Source

[Wadhwa et al., SIGGRAPH 2013]

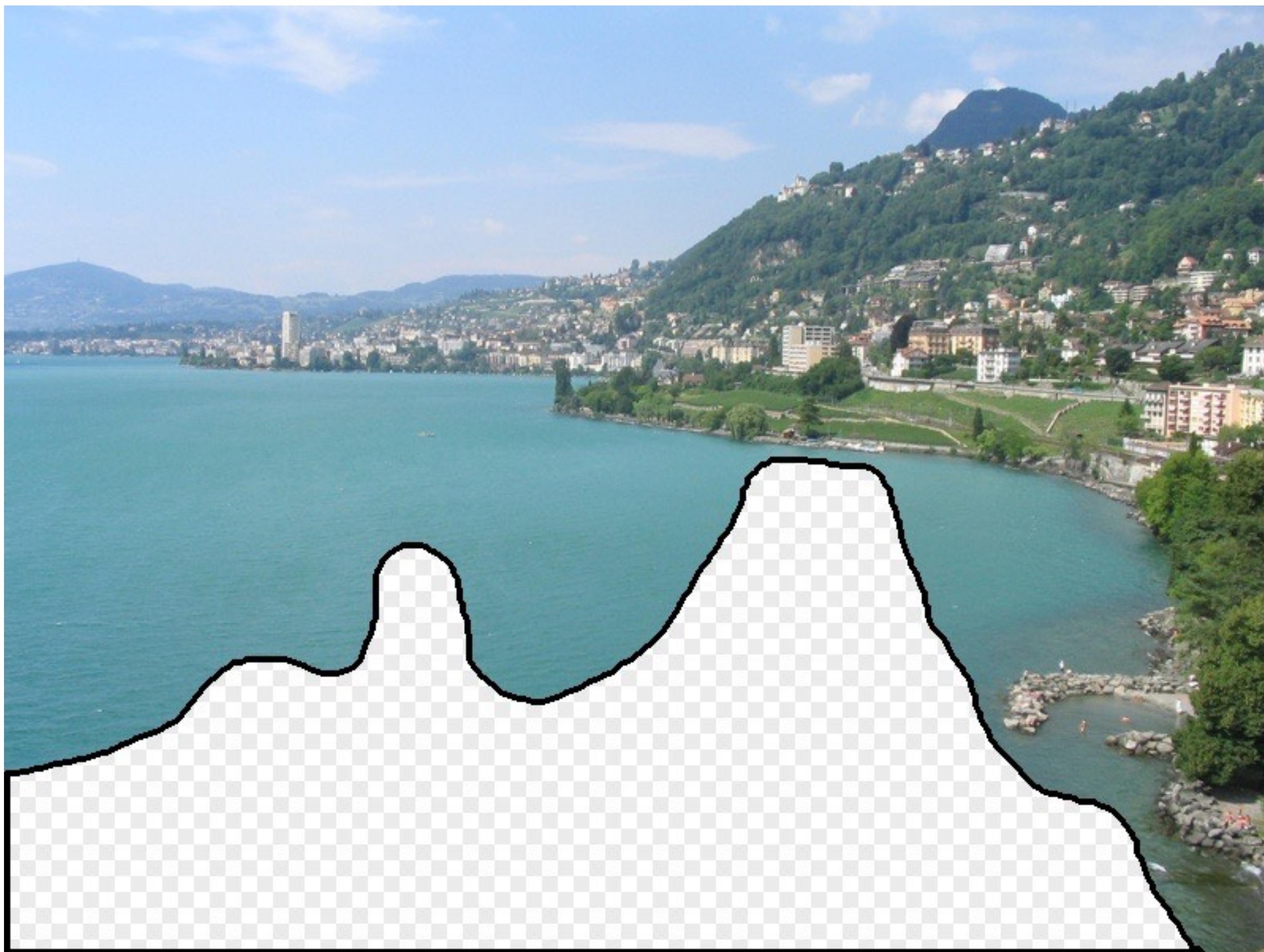
# Machine learning

# Last week: hole filling





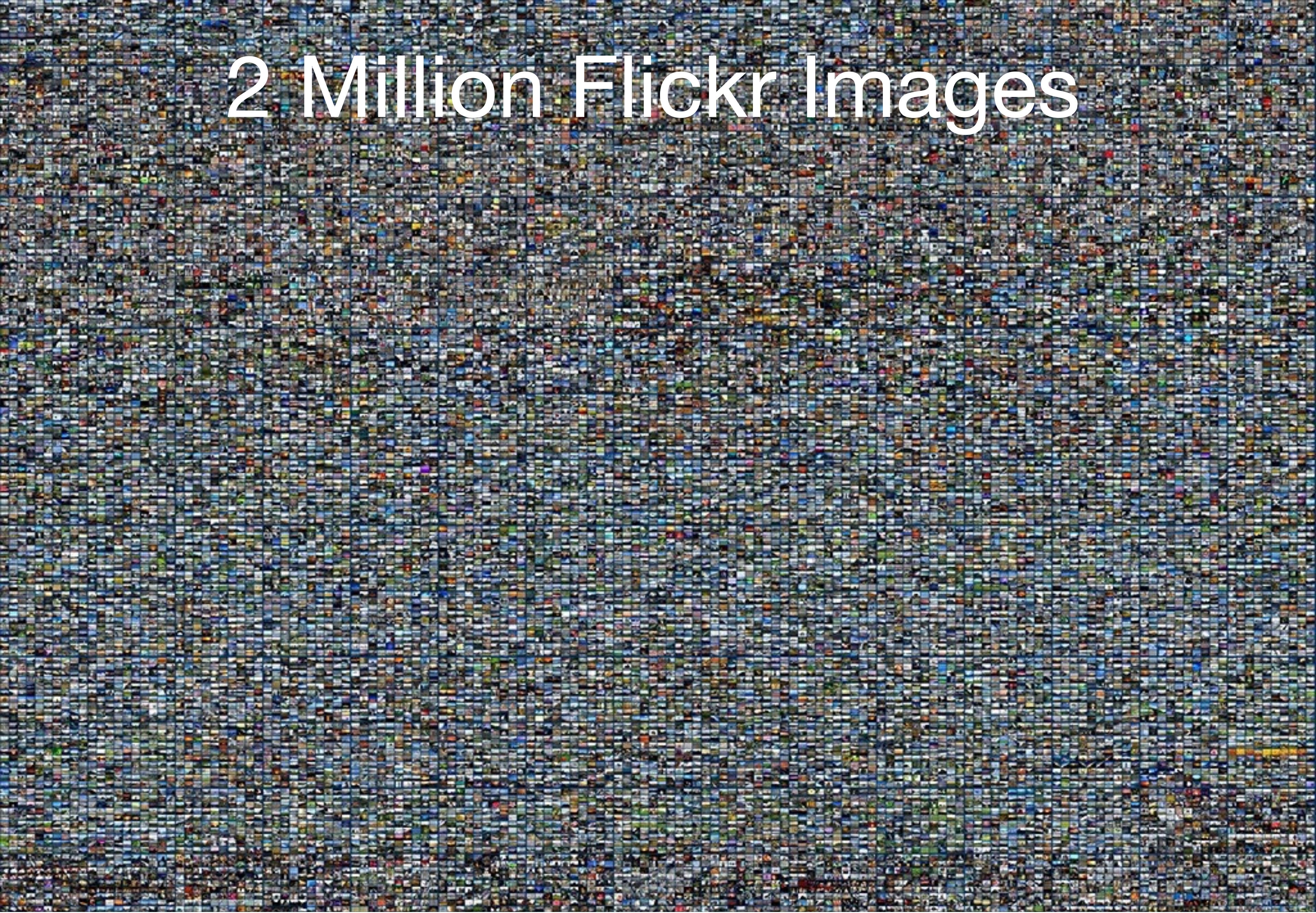
Source: A. Efros



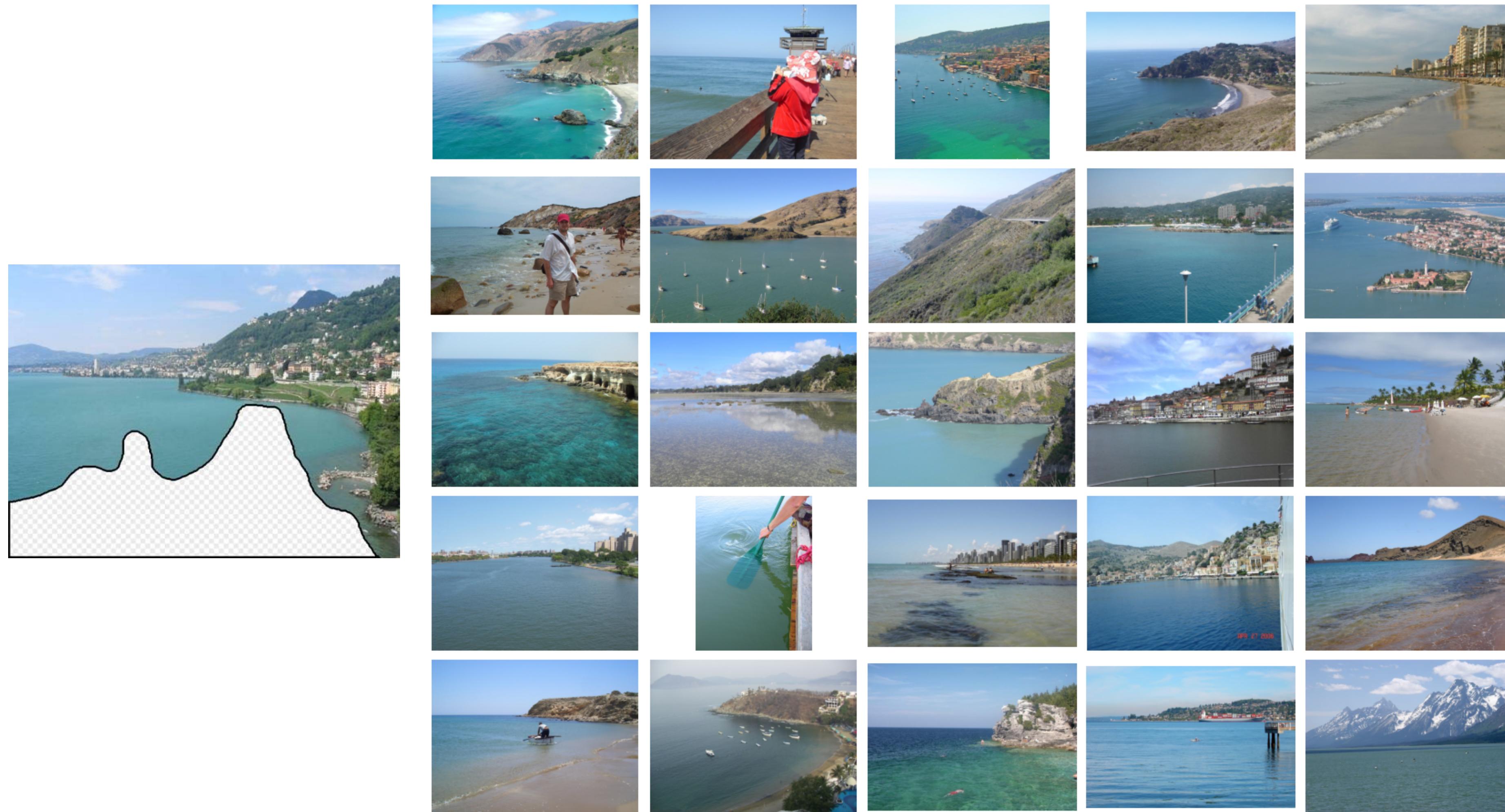
[Hays and Efros. Scene Completion Using Millions of Photographs. SIGGRAPH 2007.]

Source: A. Efros

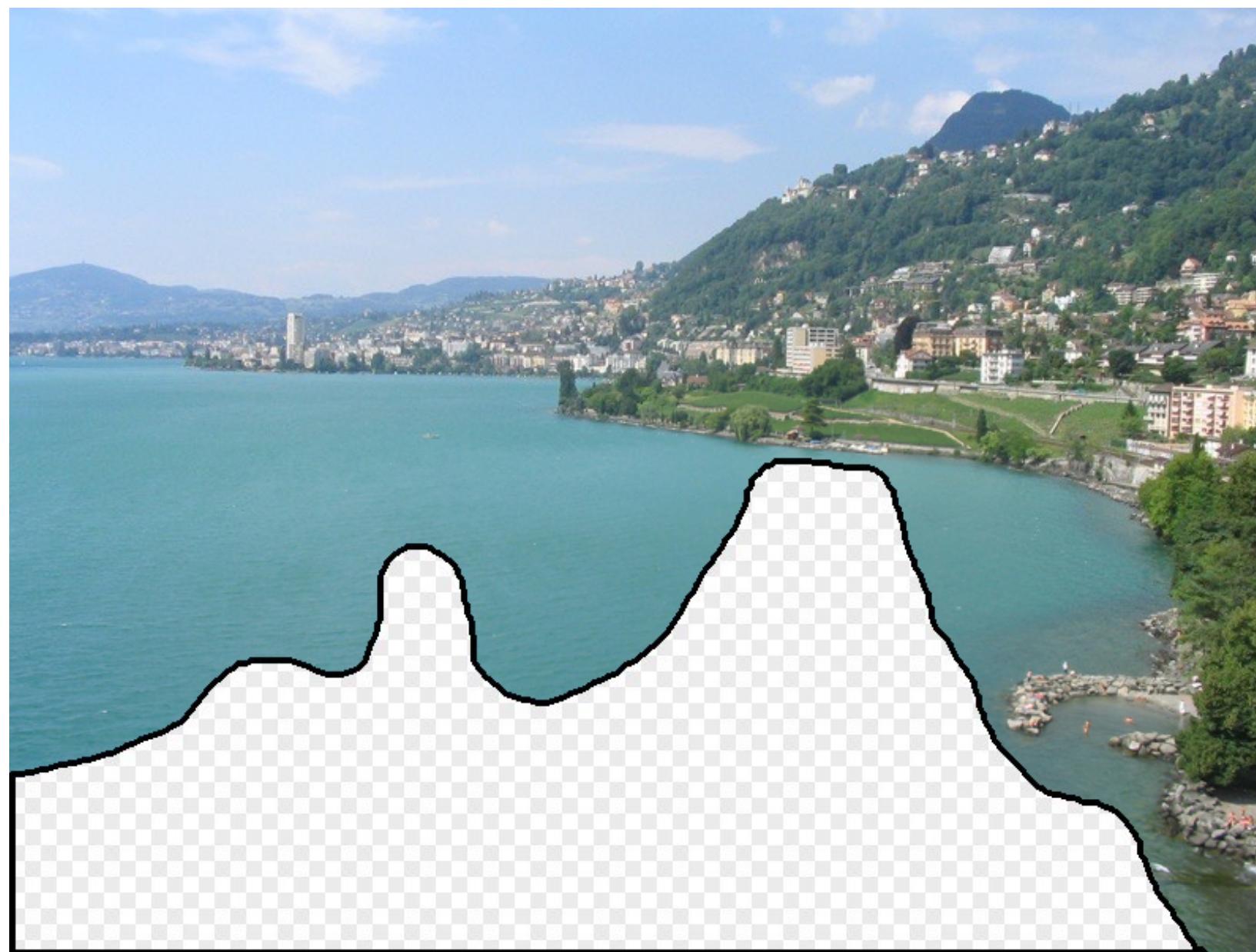
# 2 Million Flickr Images



Source: A. Efros



... 200 total



Source: A. Efros



Source: A. Efros



Source: A. Efros



Source: A. Efros



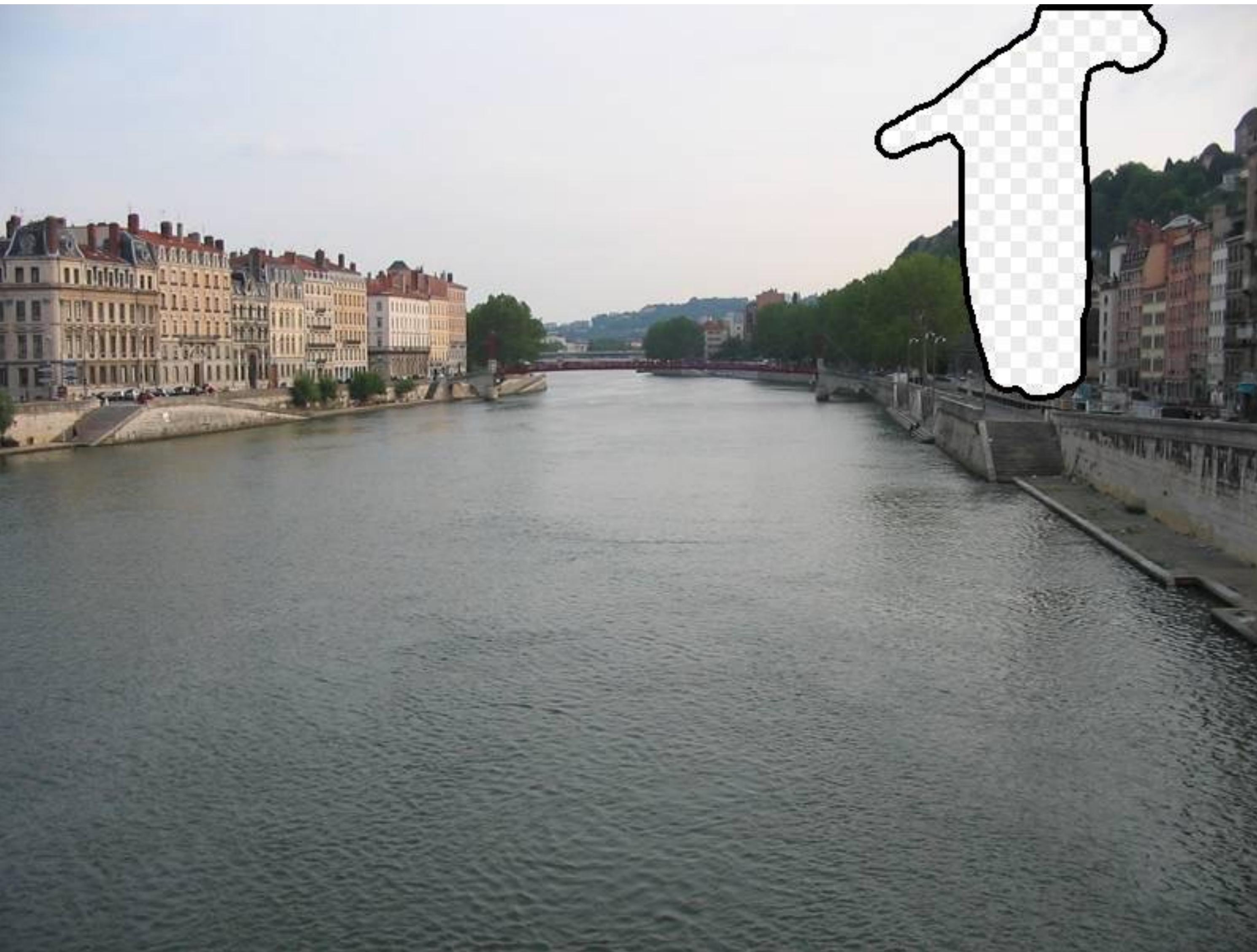
Source: A. Efros



Source: A. Efros



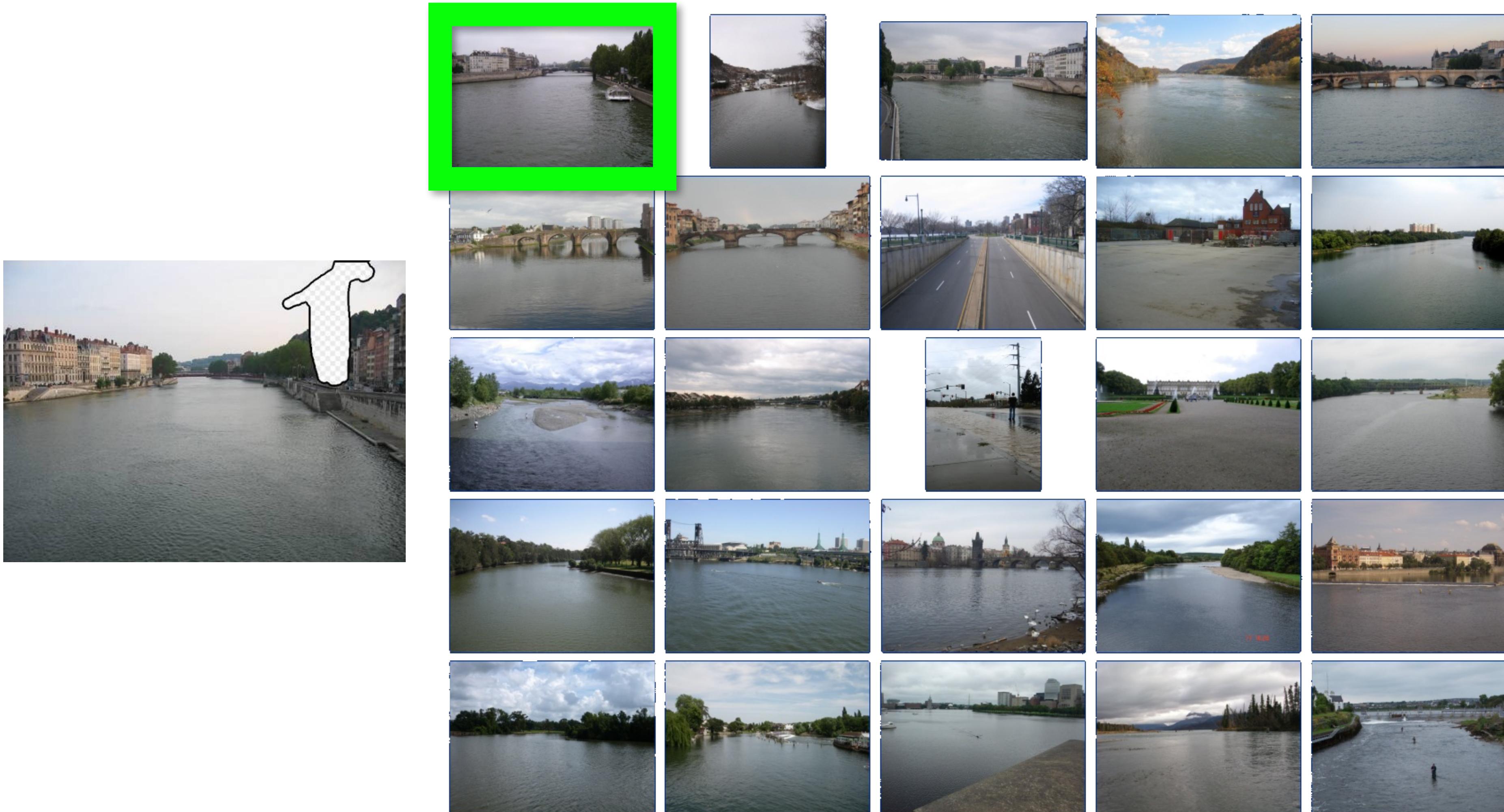
Source: A. Efros



Source: A. Efros

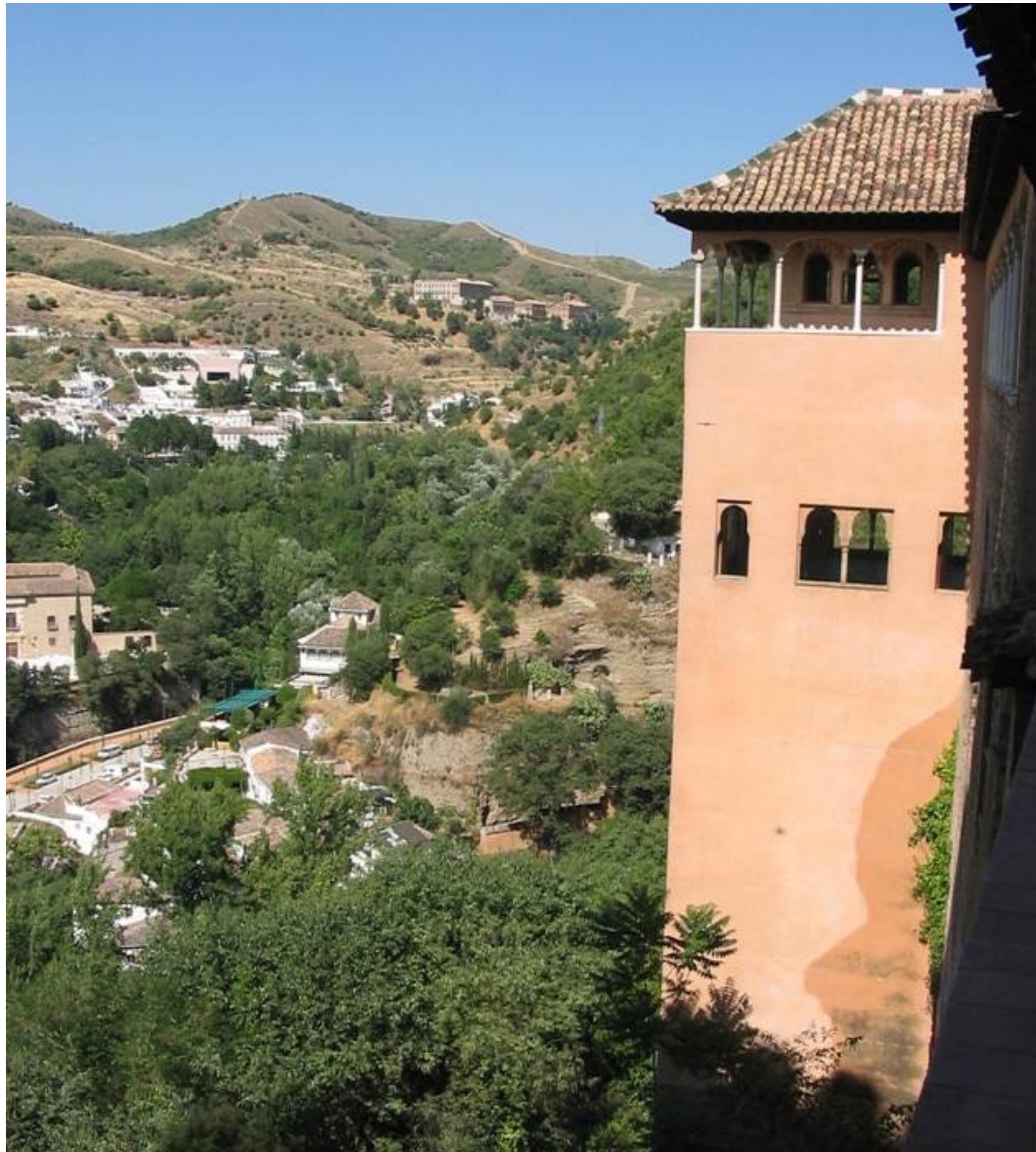


Source: A. Efros



... 200 scene matches





Source: A. Efros

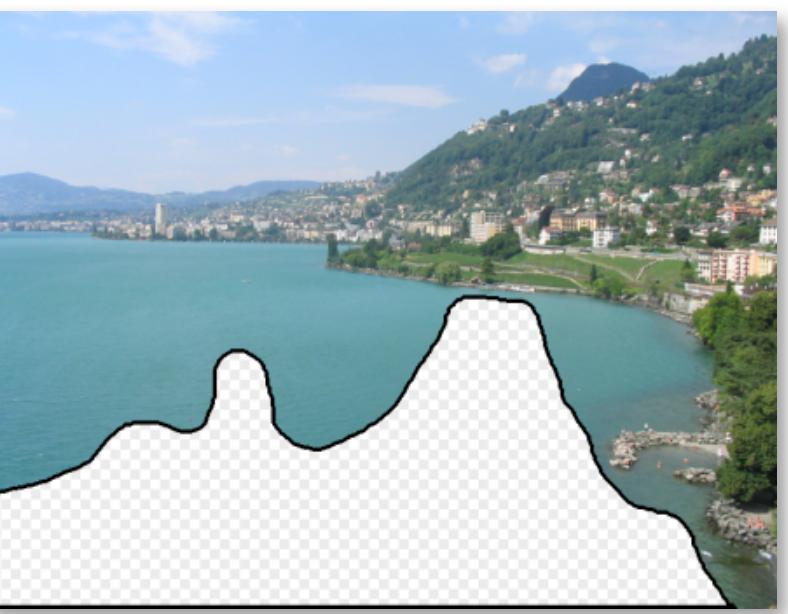


Source: A. Efros

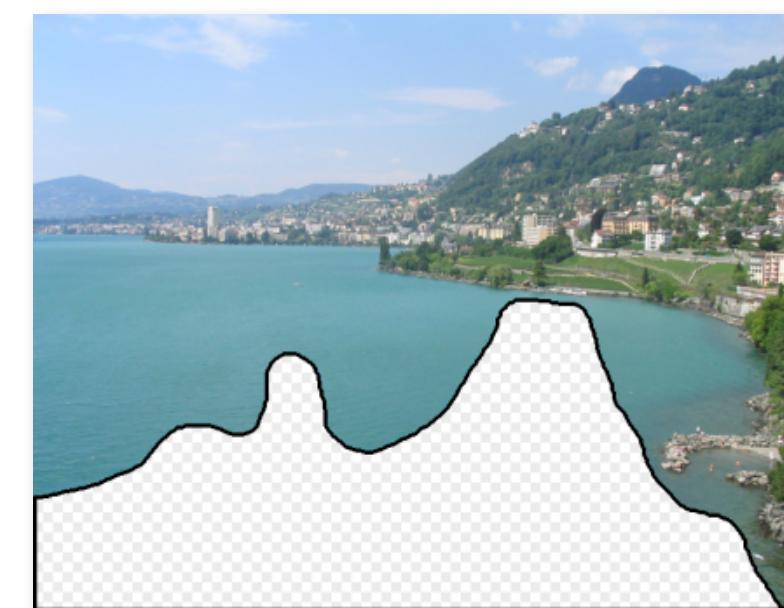


Source: A. Efros

# Why does it work?

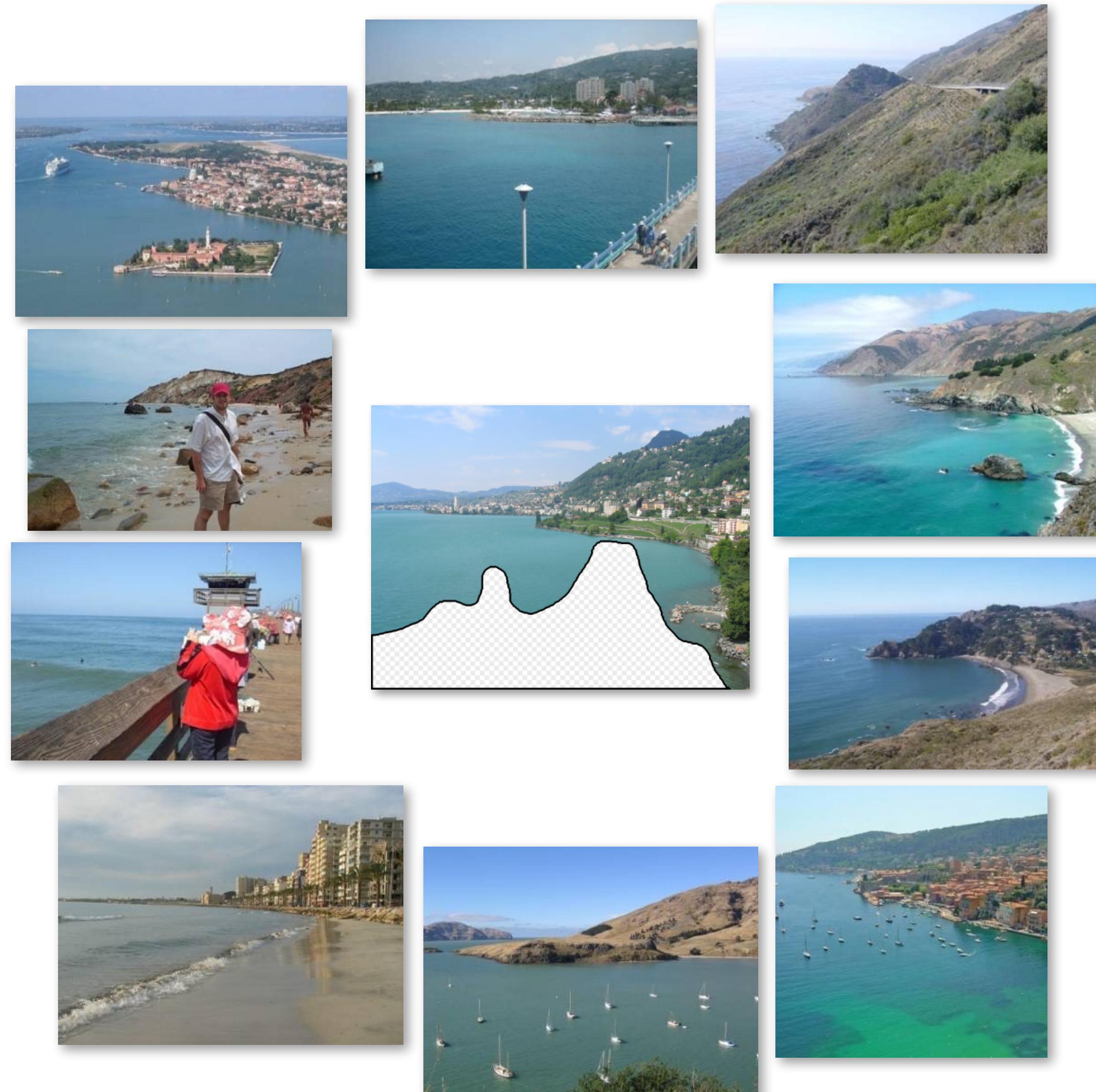


Source: A. Efros



Nearest neighbors from a collection of 20 thousand images

Source: A. Efros



Nearest neighbors from a  
collection of 2 million images

Source: A. Efros

# “Unreasonable Effectiveness of Data”

[Halevy, Norvig, Pereira 2009]

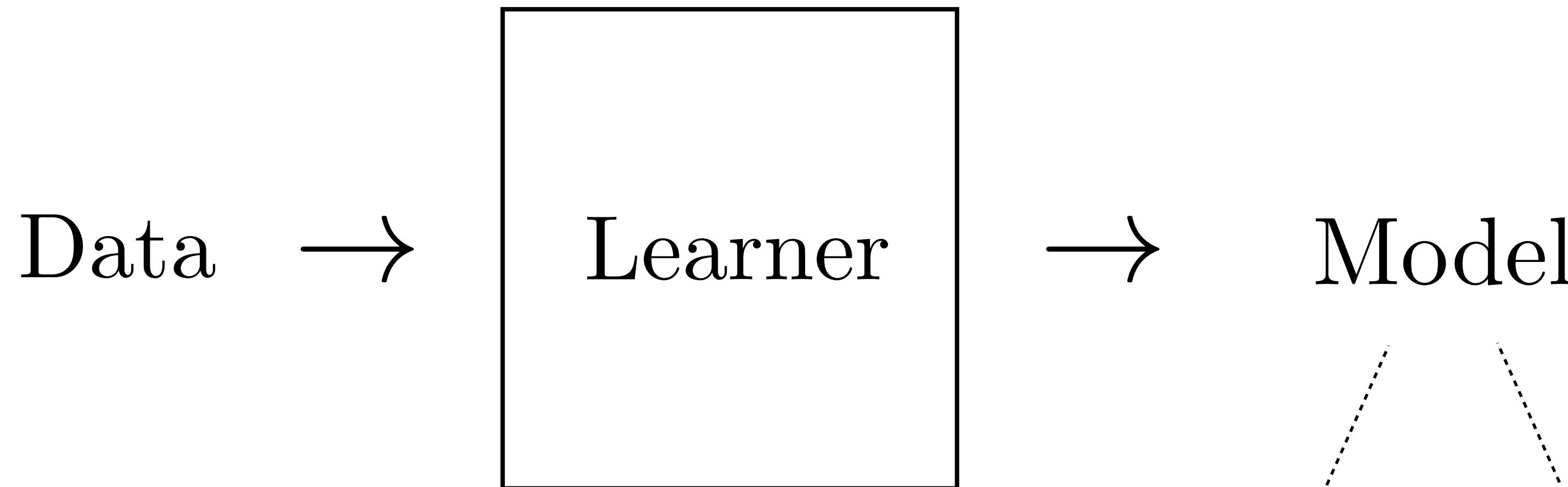
Parts of our world can be explained by elegant  
mathematics  
physics, chemistry, astronomy, etc.

But much cannot  
psychology, economics, genetics, etc.

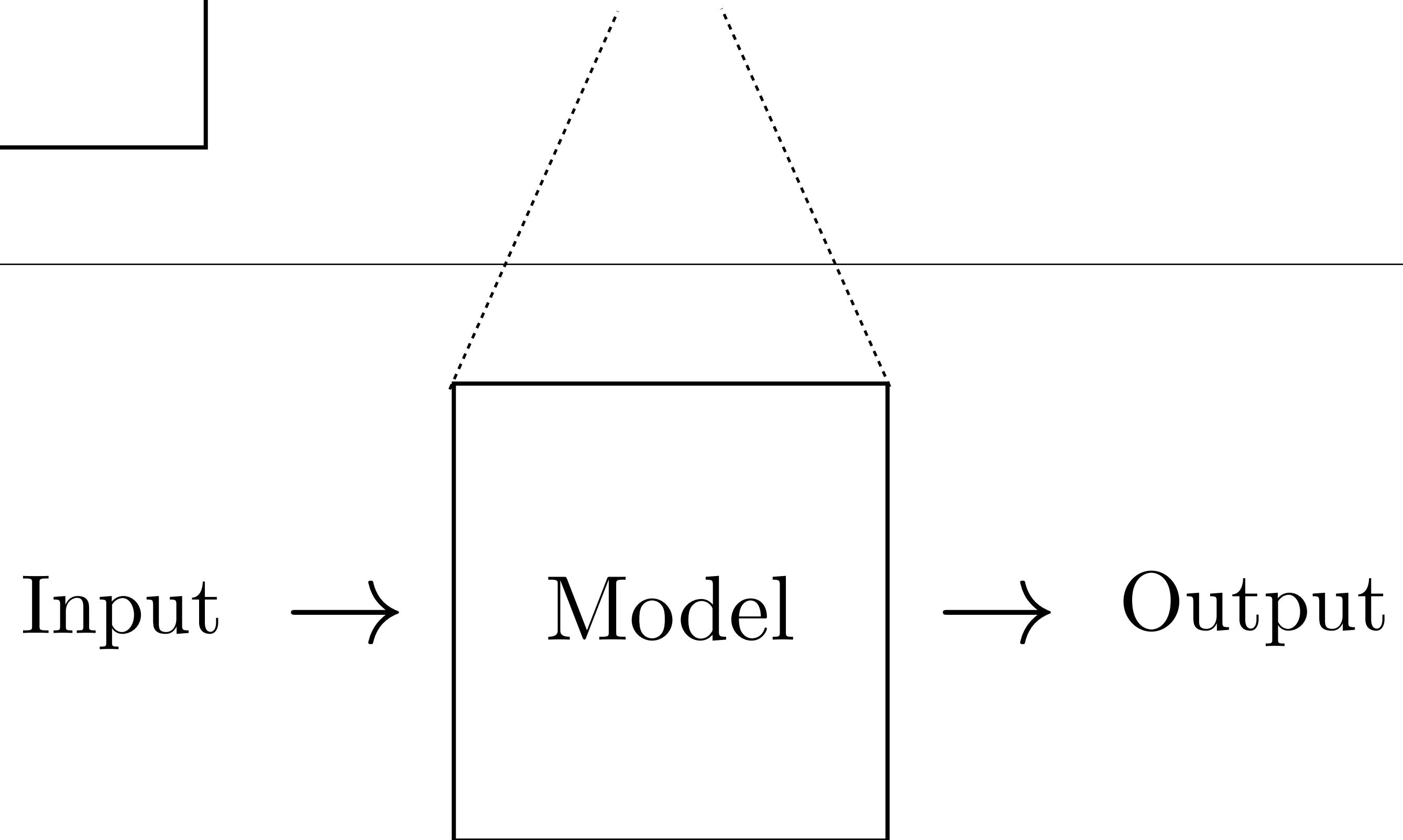
Enter The Data!

“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



# What does $\star$ do?

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

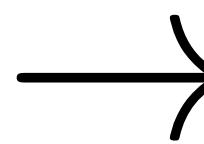
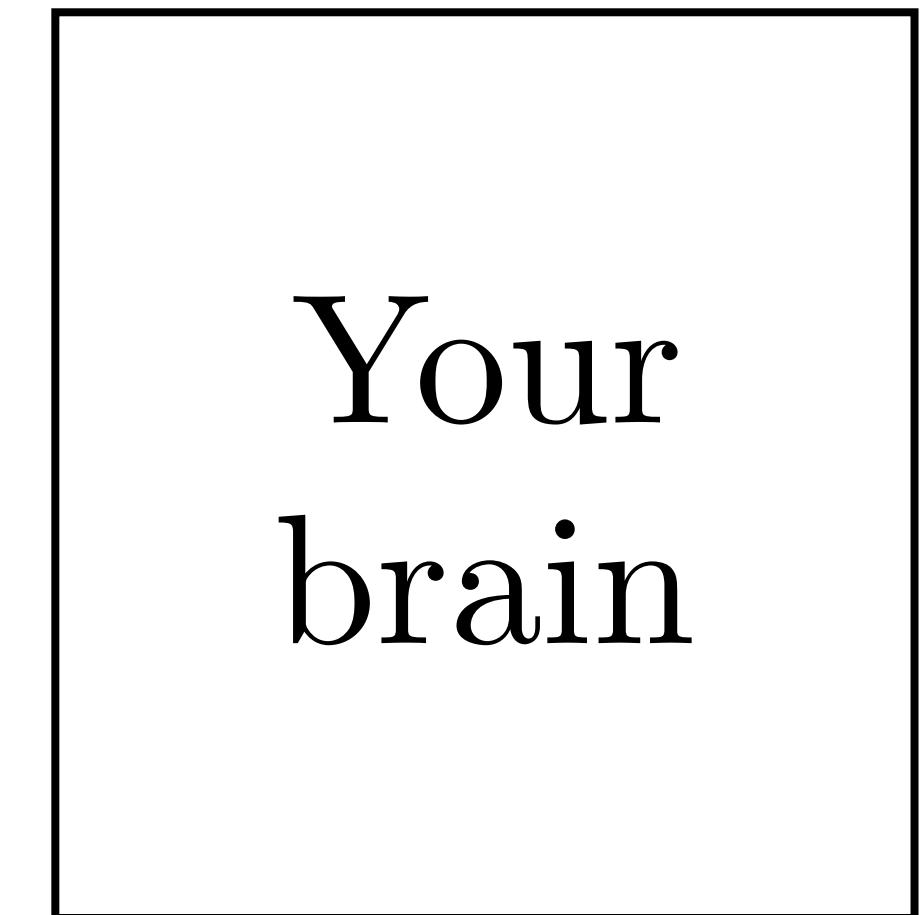
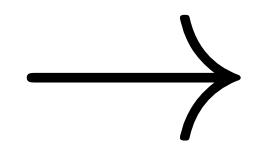
## Training

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$



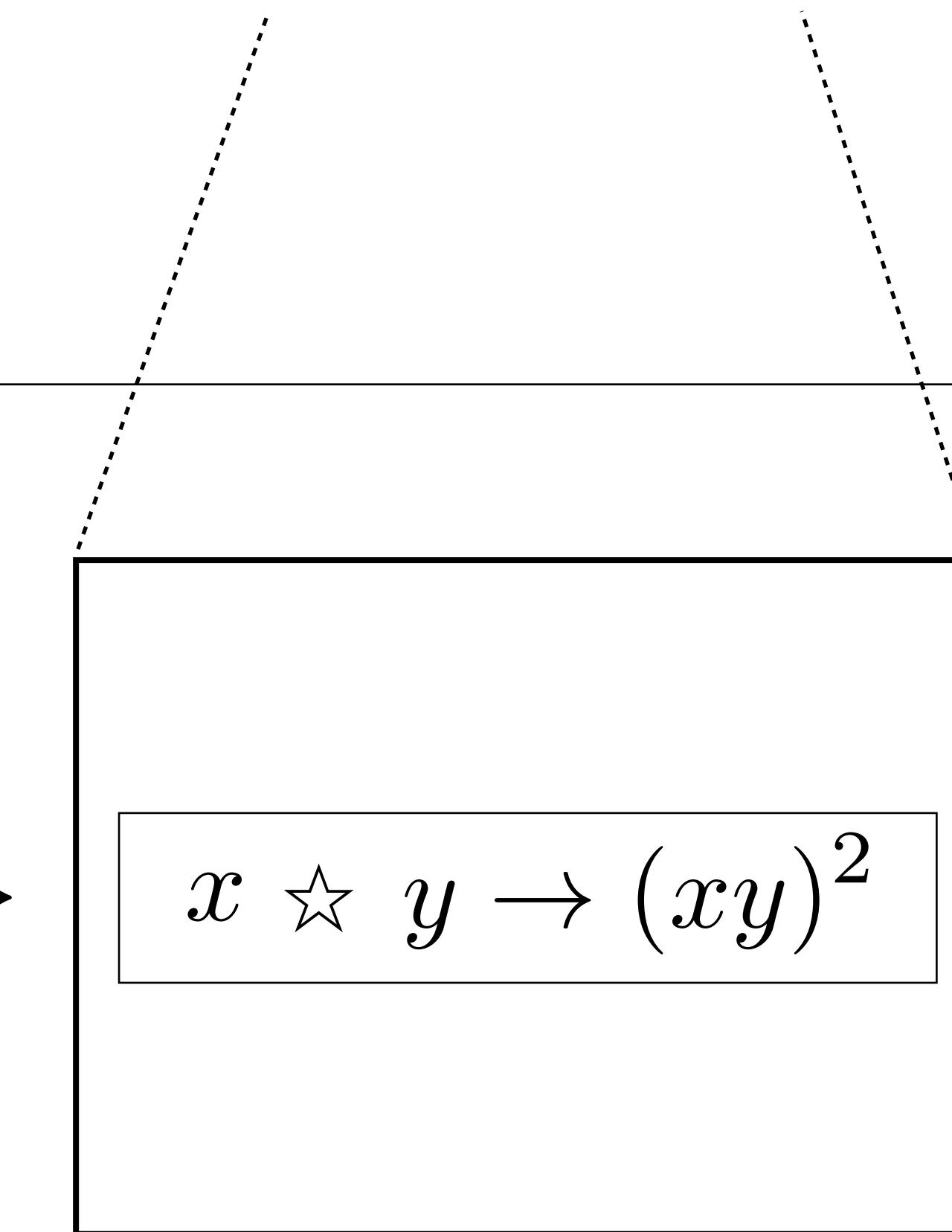
$$x \star y \rightarrow (xy)^2$$

## Testing

$$3 \star 5 \rightarrow$$

$$x \star y \rightarrow (xy)^2$$

$$\rightarrow 225$$

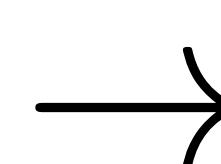
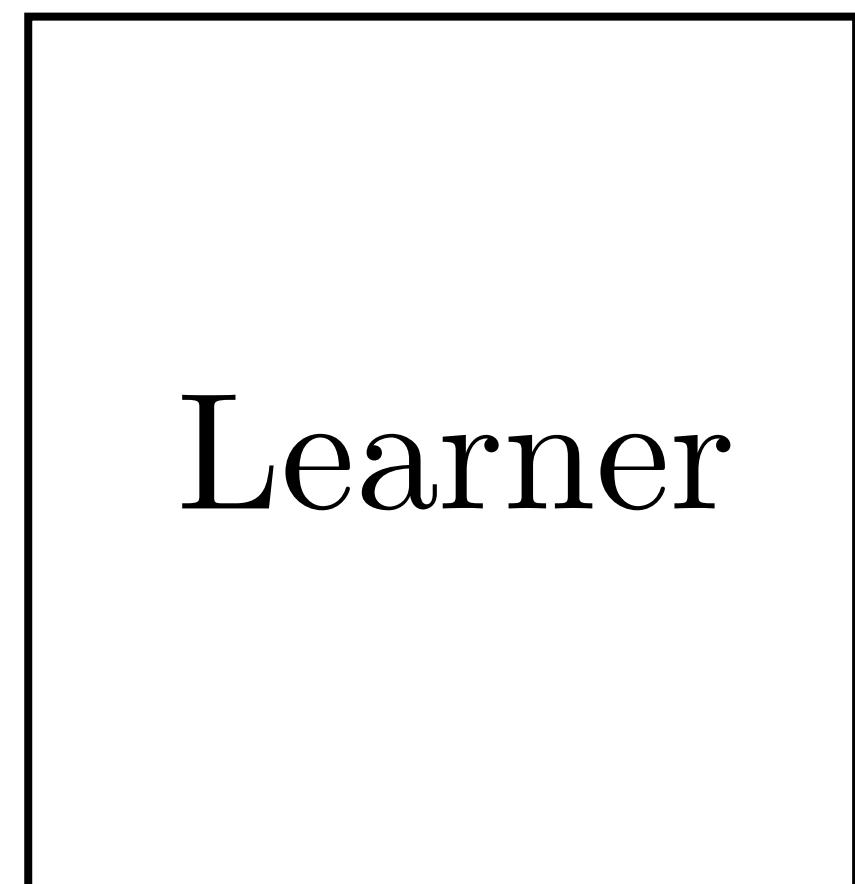


# Learning from examples

(aka **supervised learning**)

Training data

```
{input:[2, 3], output:36}  
{input:[7, 1], output:49}  
{input:[5, 2], output:100}  
{input:[2, 2], output:16}
```


$$f$$

# Learning from examples

(aka **supervised learning**)

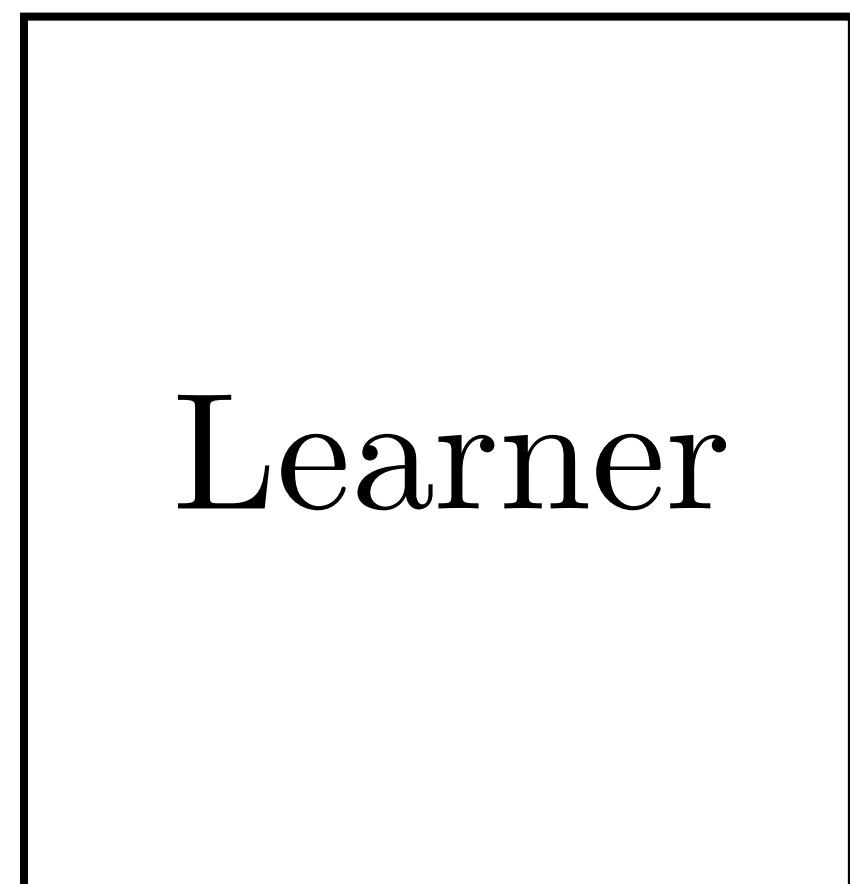
Training data

$$\{x_1, y_1\}$$

$$\{x_2, y_2\} \rightarrow$$

$$\{x_3, y_3\}$$

...

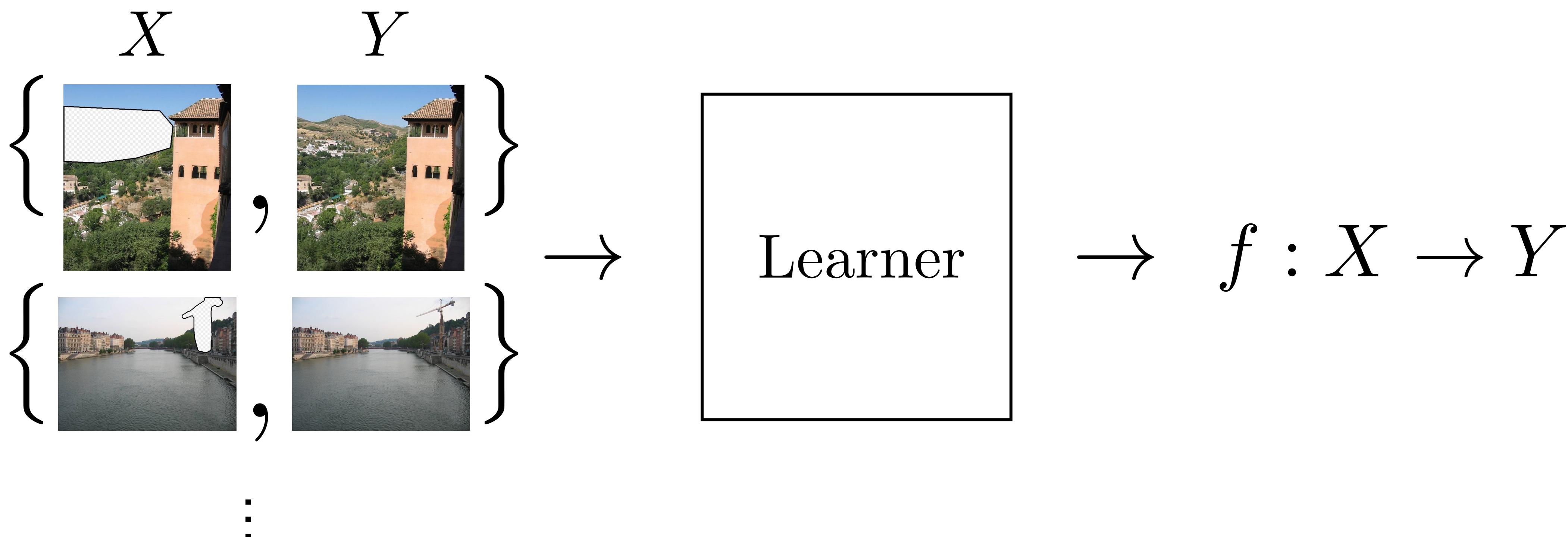


$$\rightarrow f : X \rightarrow Y$$

# Learning from examples

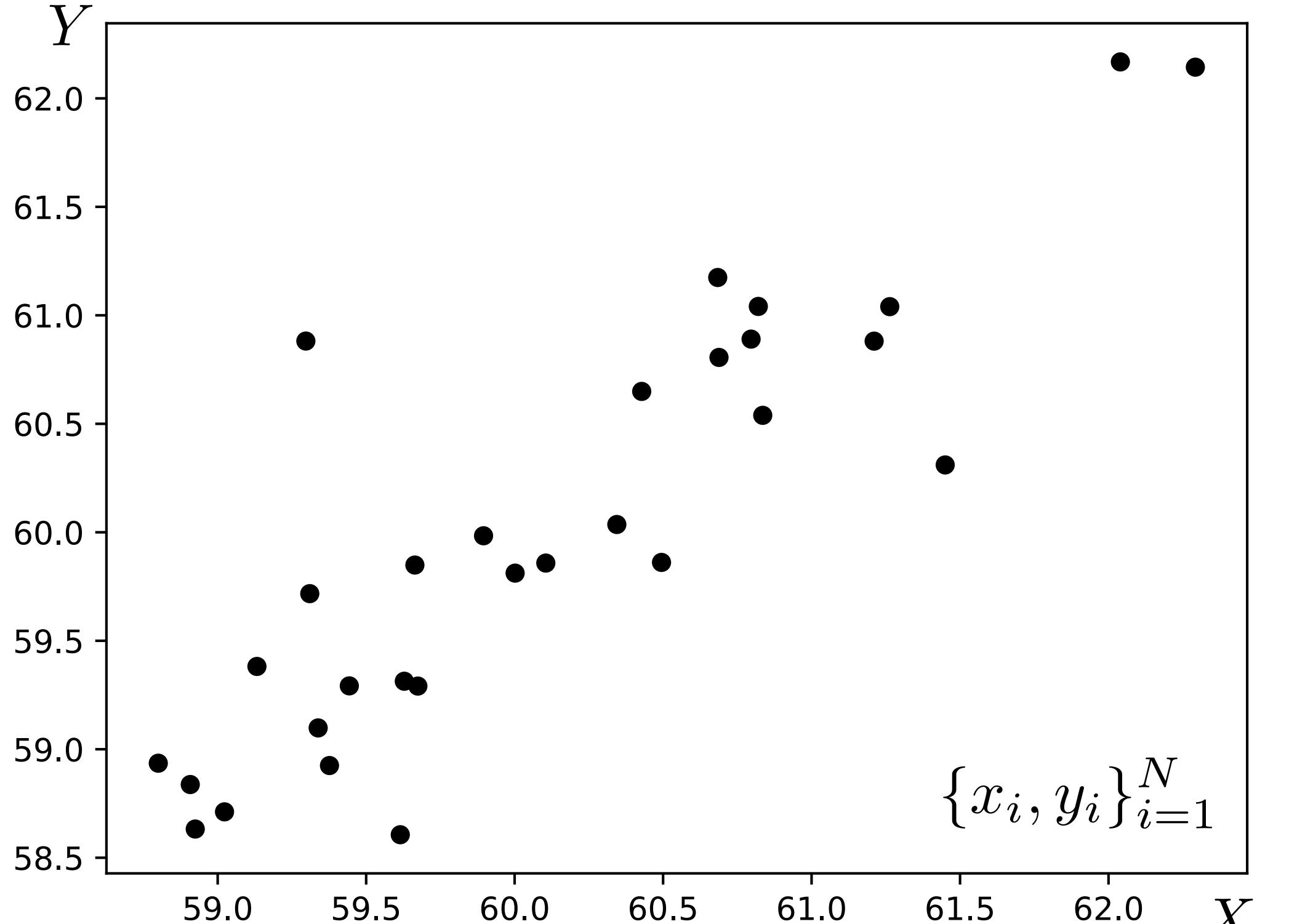
(aka **supervised learning**)

Training data

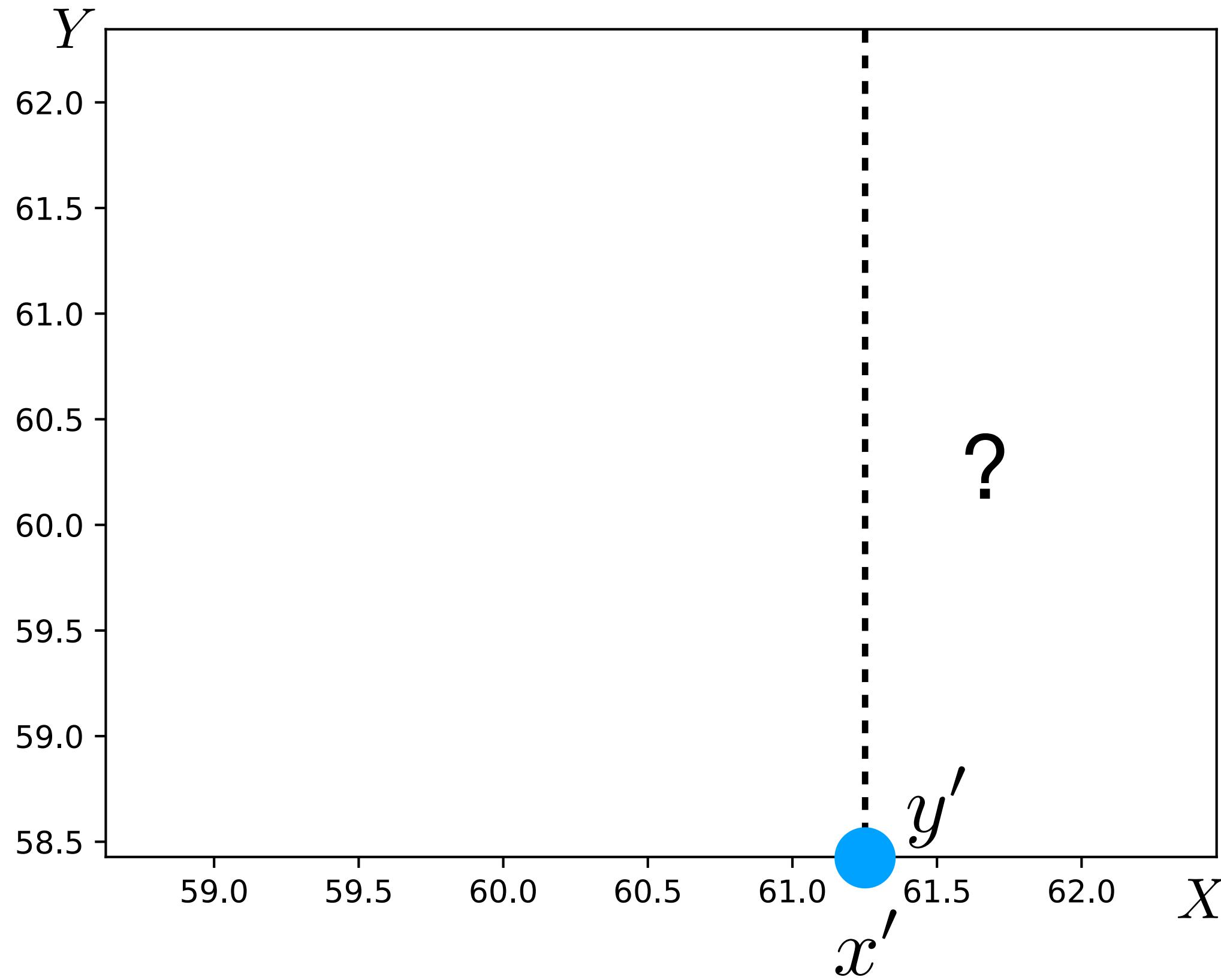


# Case study #1: Linear least squares

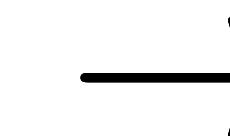
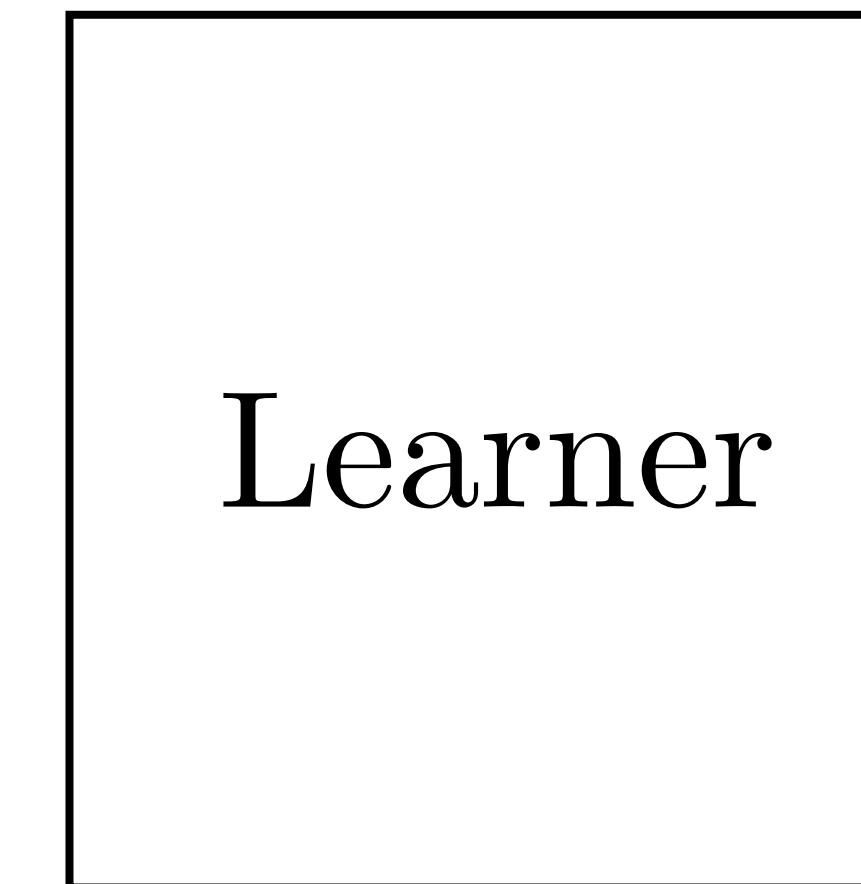
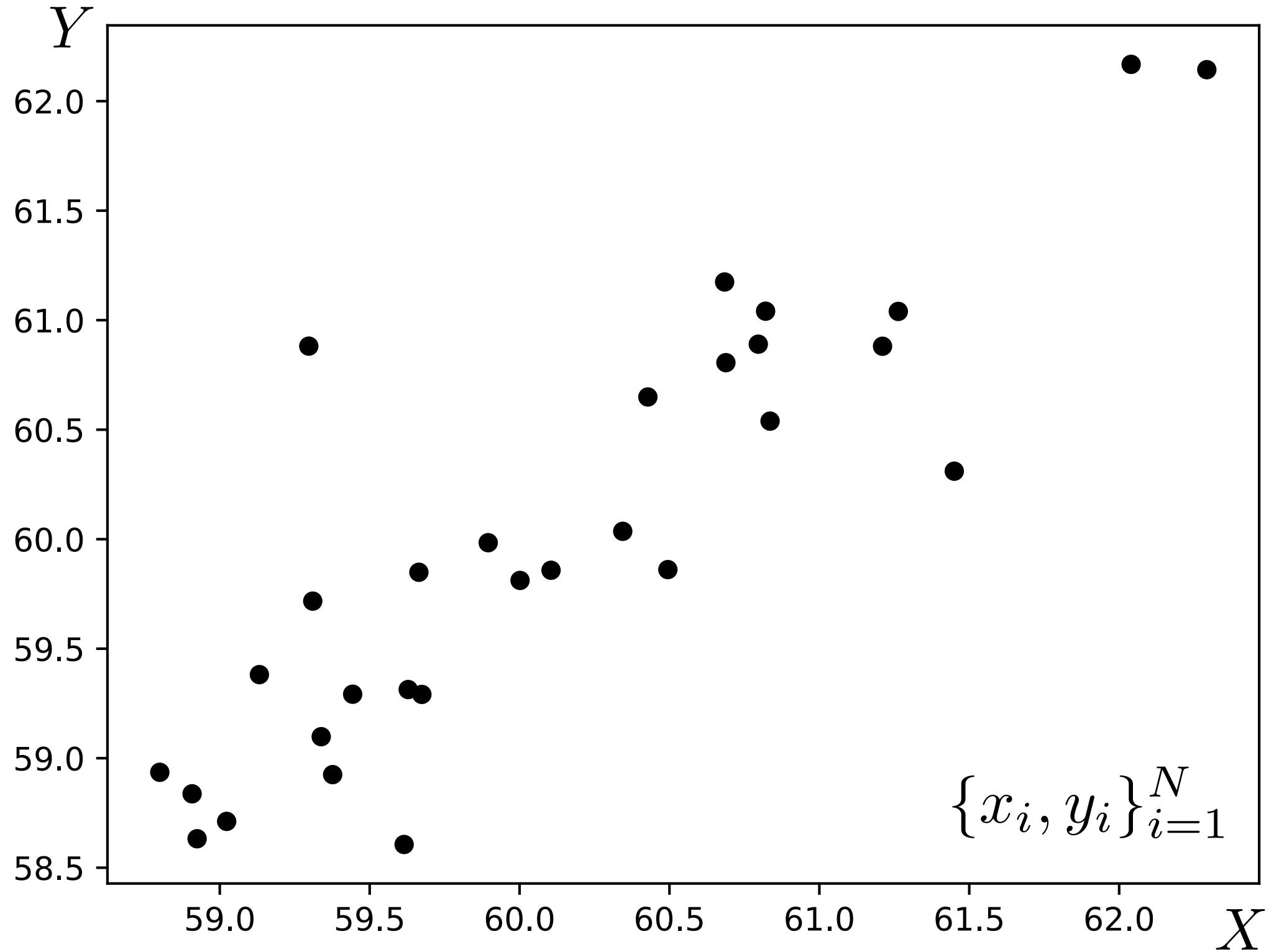
# Training data



# Test query



## Training data

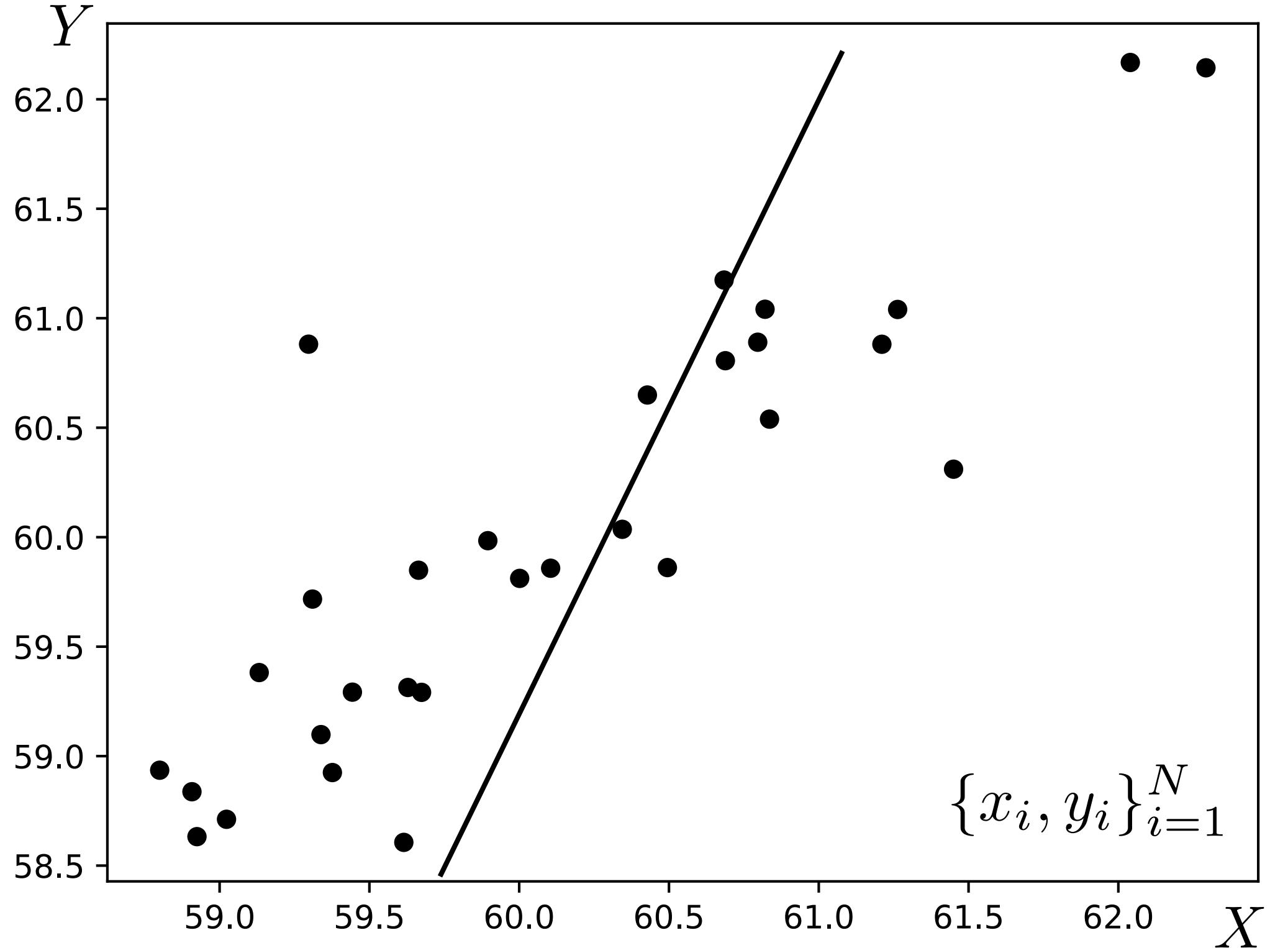


$$f_{\theta}(x) = \theta_1 x + \theta_0$$

## Hypothesis space

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

# Training data

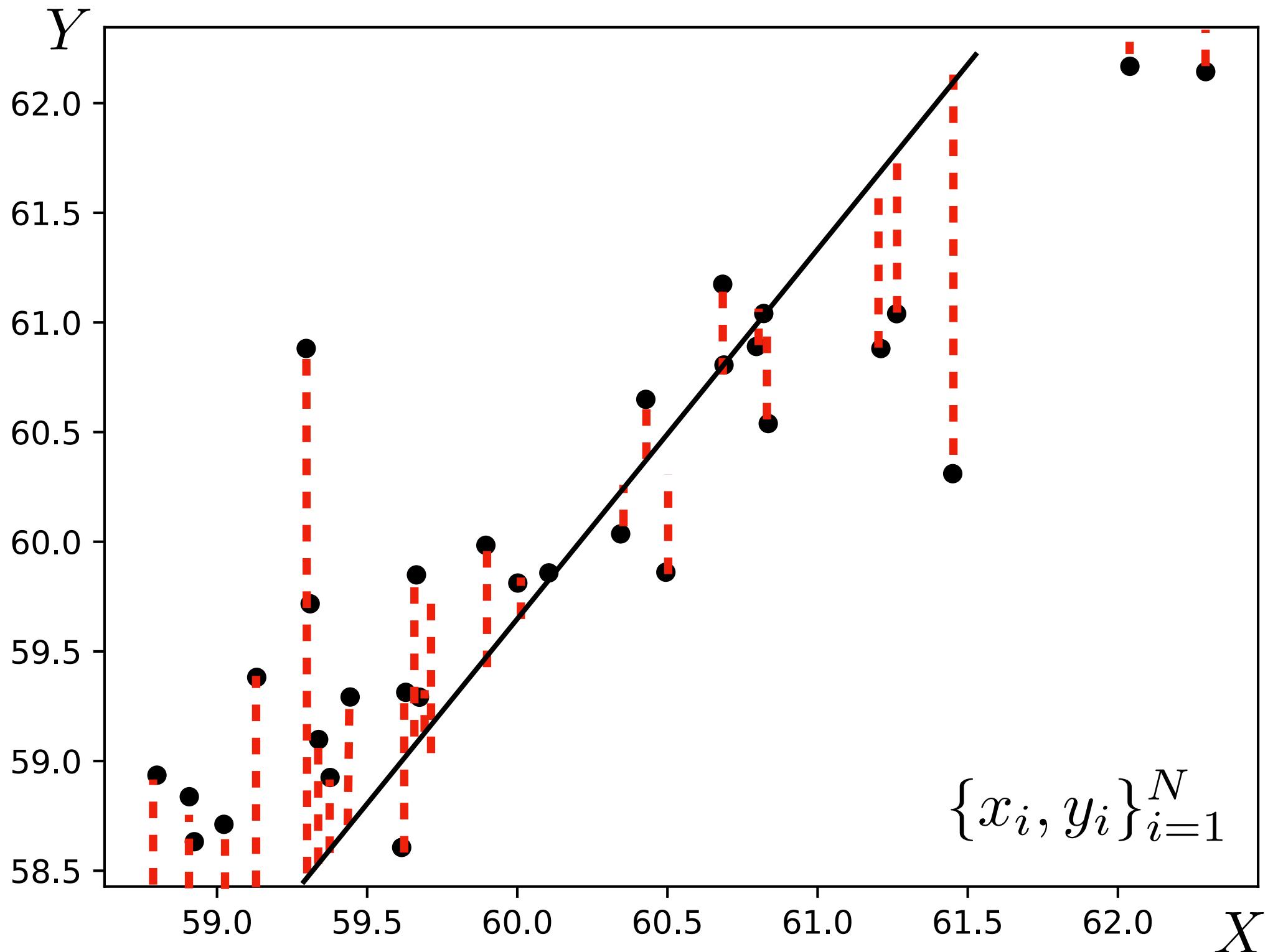


Search for the **parameters**,  $\theta = \{\theta_0, \theta_1\}$ , that best fit the data.

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

Best fit in what sense?

# Training data



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

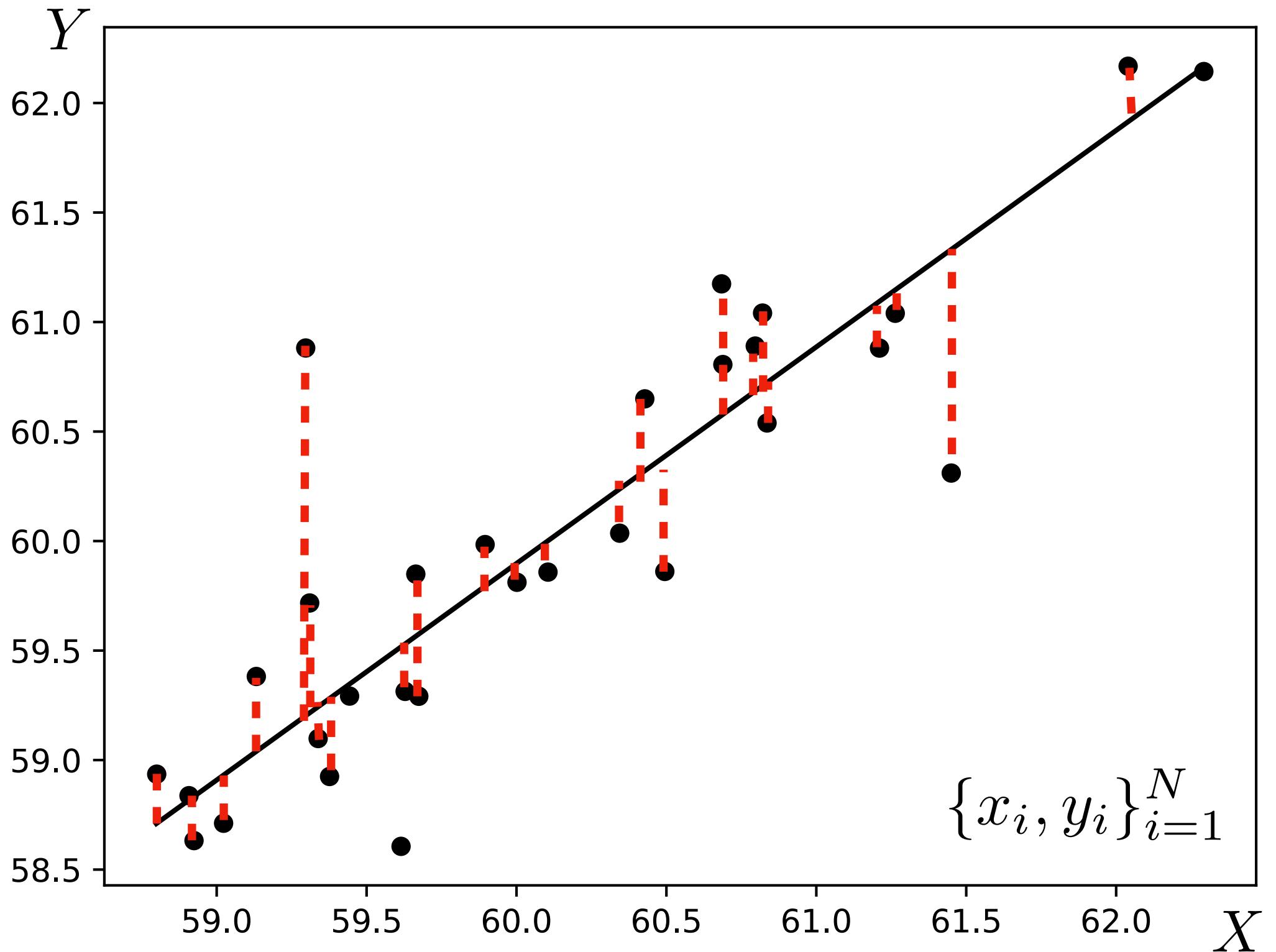
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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.

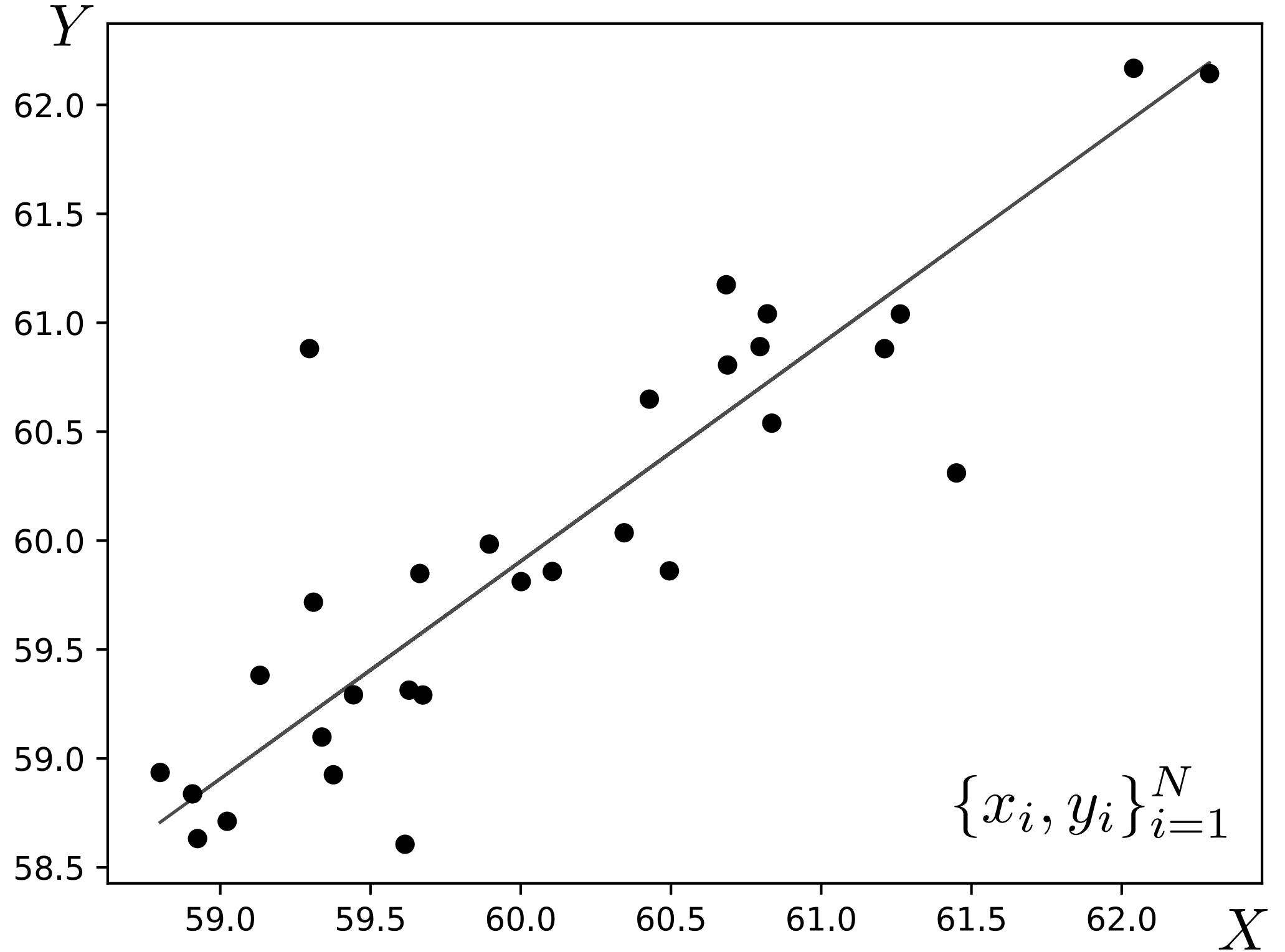
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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

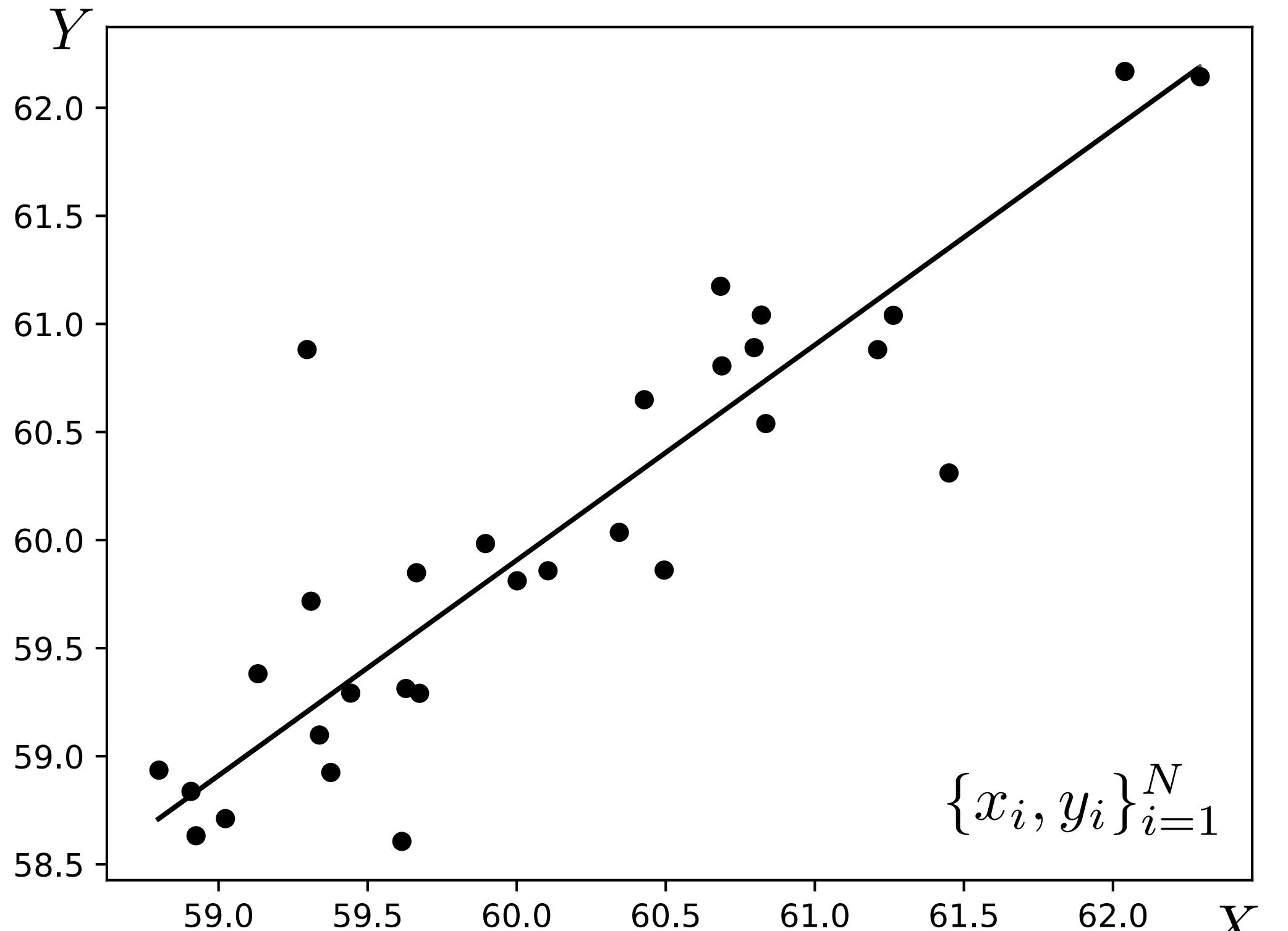


**Complete learning problem:**

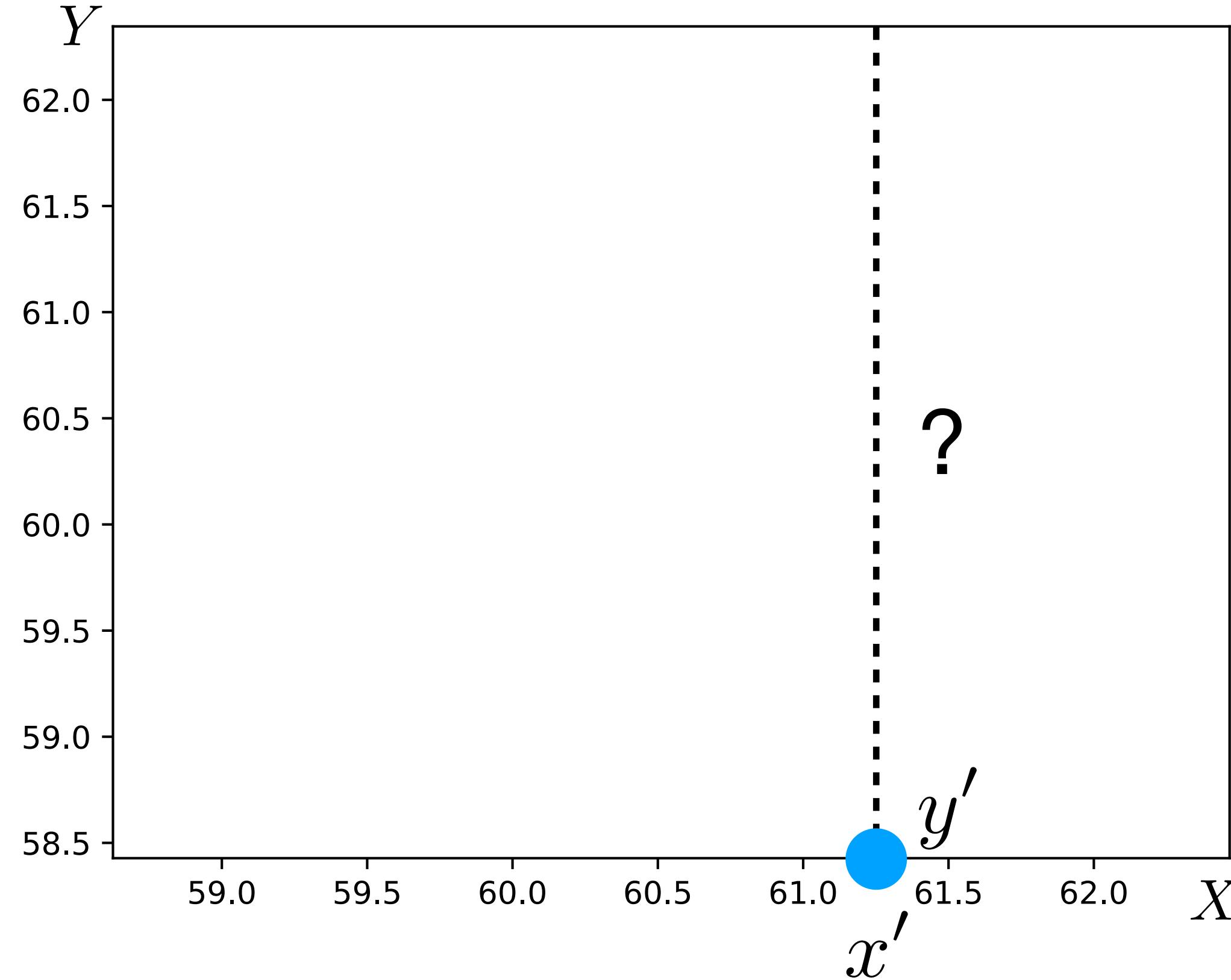
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x_i) - y_i)^2$$

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

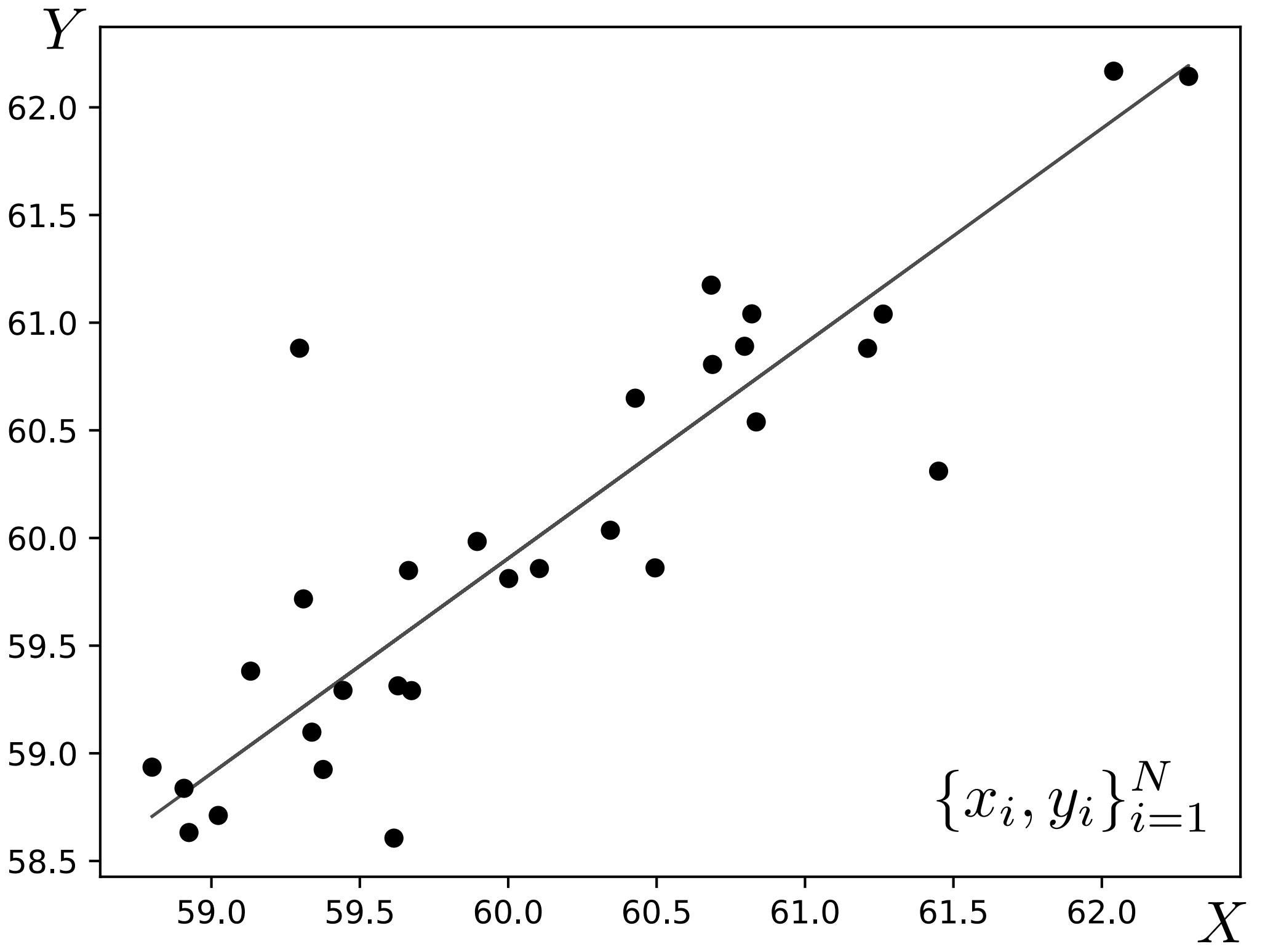
# Training data



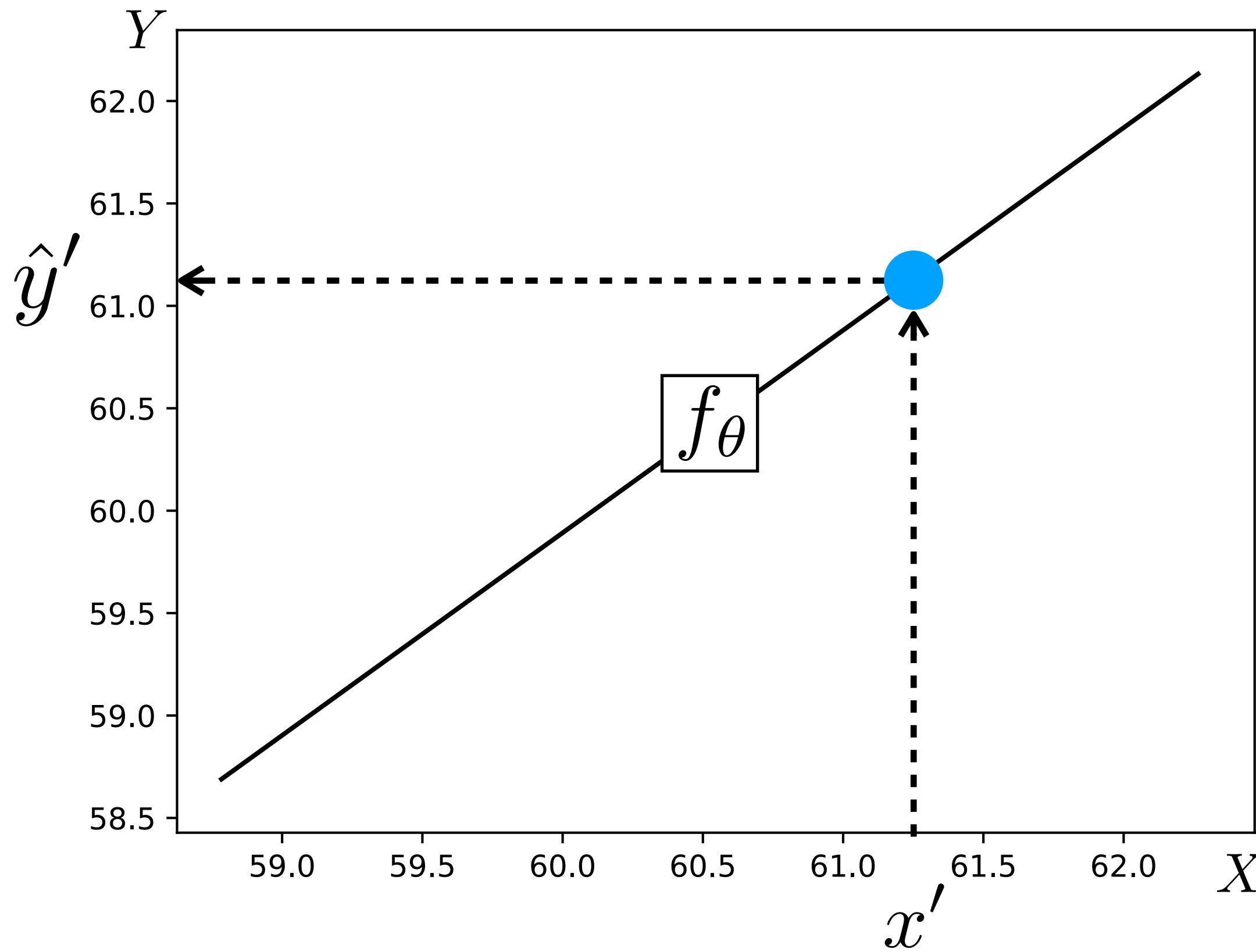
# Test query



# Training data

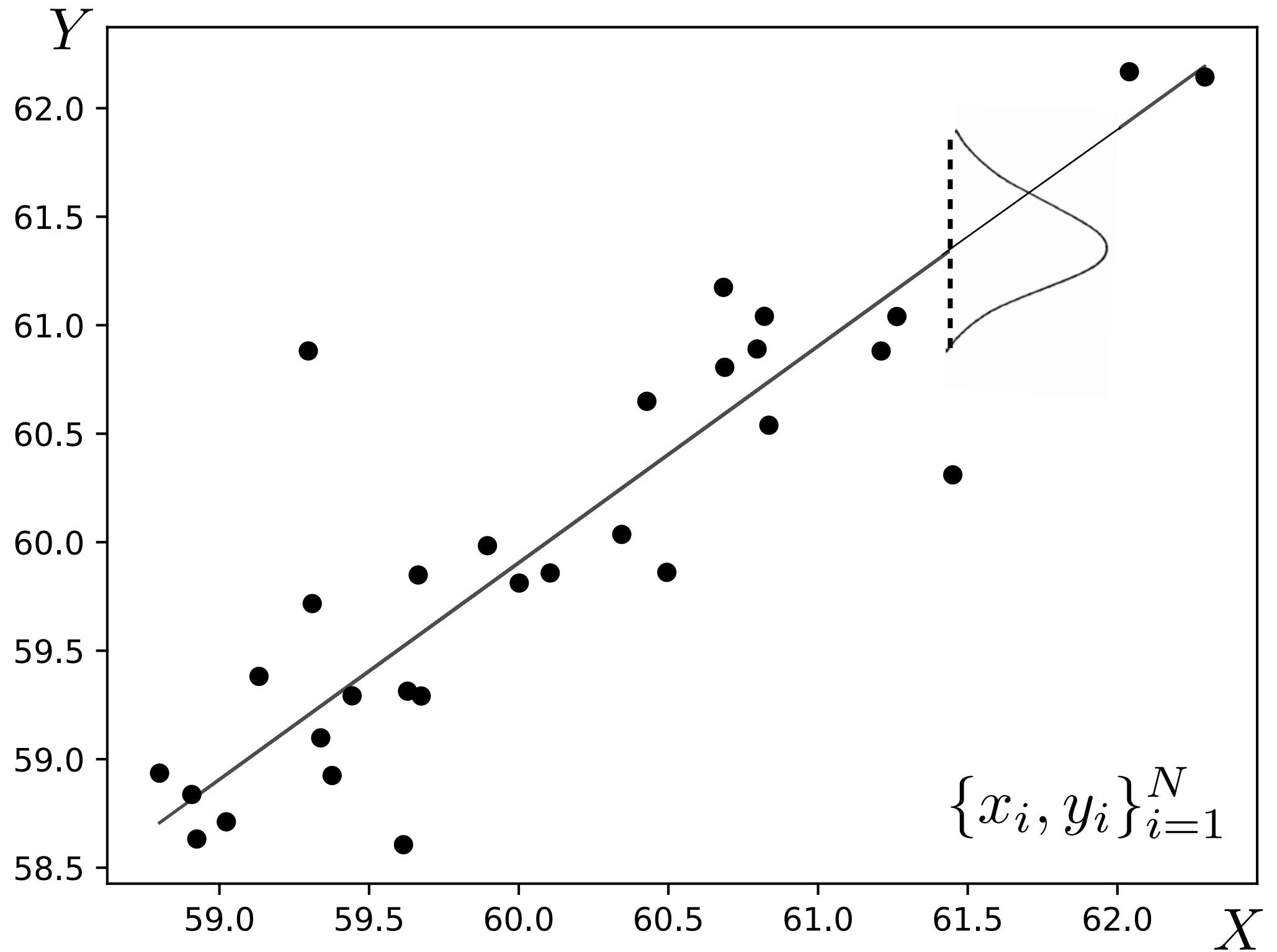


# Test query



$$x' \dashrightarrow [f_\theta] \dashrightarrow \hat{y}'$$

# Why is squared error a good objective?



**Follows from two *modeling assumptions***

1. Observed Y's are corrupted by Gaussian noise:

$$Y \approx f_\theta(X)$$

$$\epsilon = Y - f_\theta(X)$$

$$\epsilon \sim \mathcal{N}(0, \sigma)$$

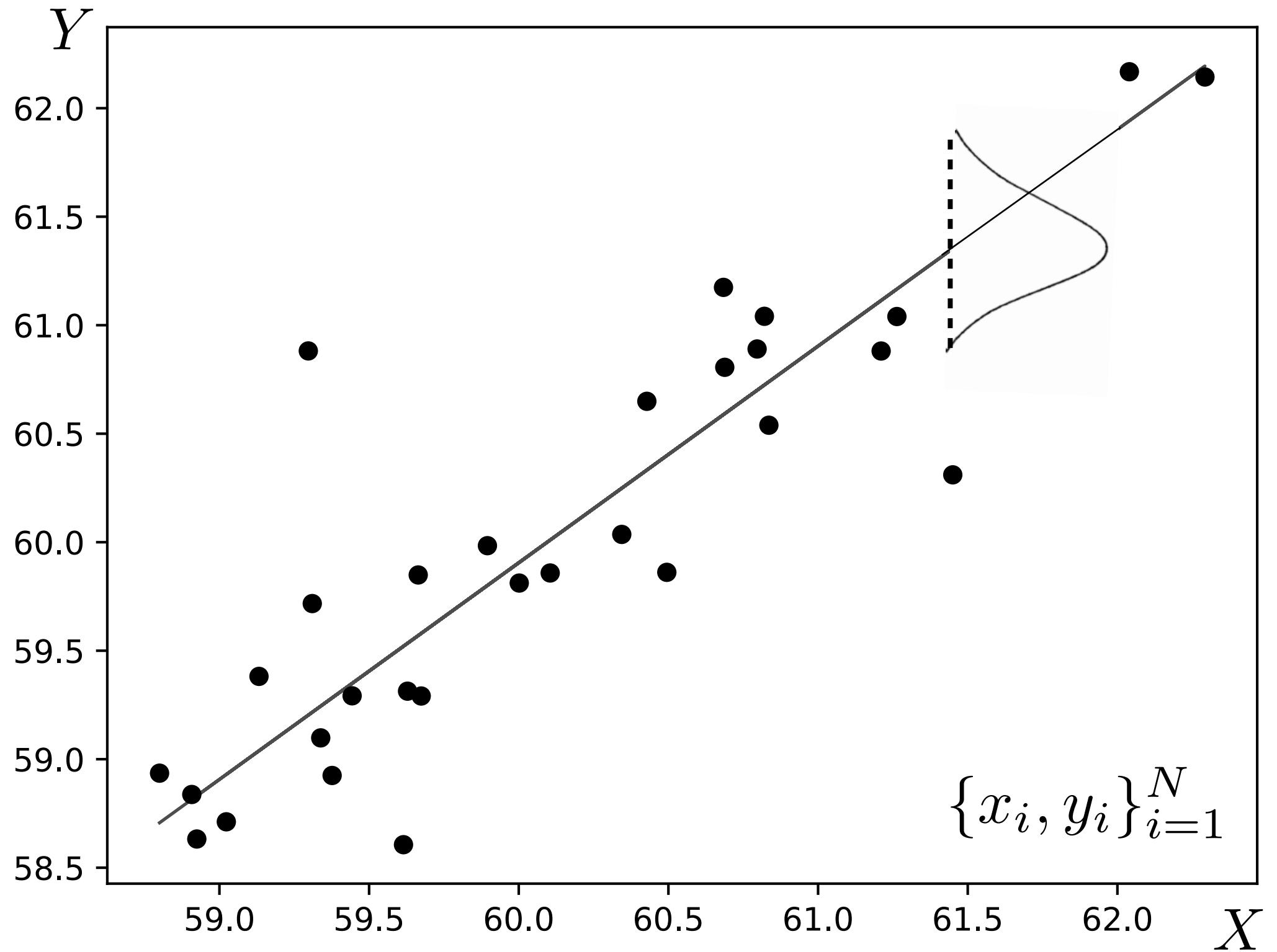
$$Y = f_\theta(X) + \epsilon$$

$$P_\theta(Y = y | X = x) \propto \exp \frac{-(y - f_\theta(x))^2}{2\sigma^2}$$

2. Training datapoints are independently and identically distributed (**iid**):

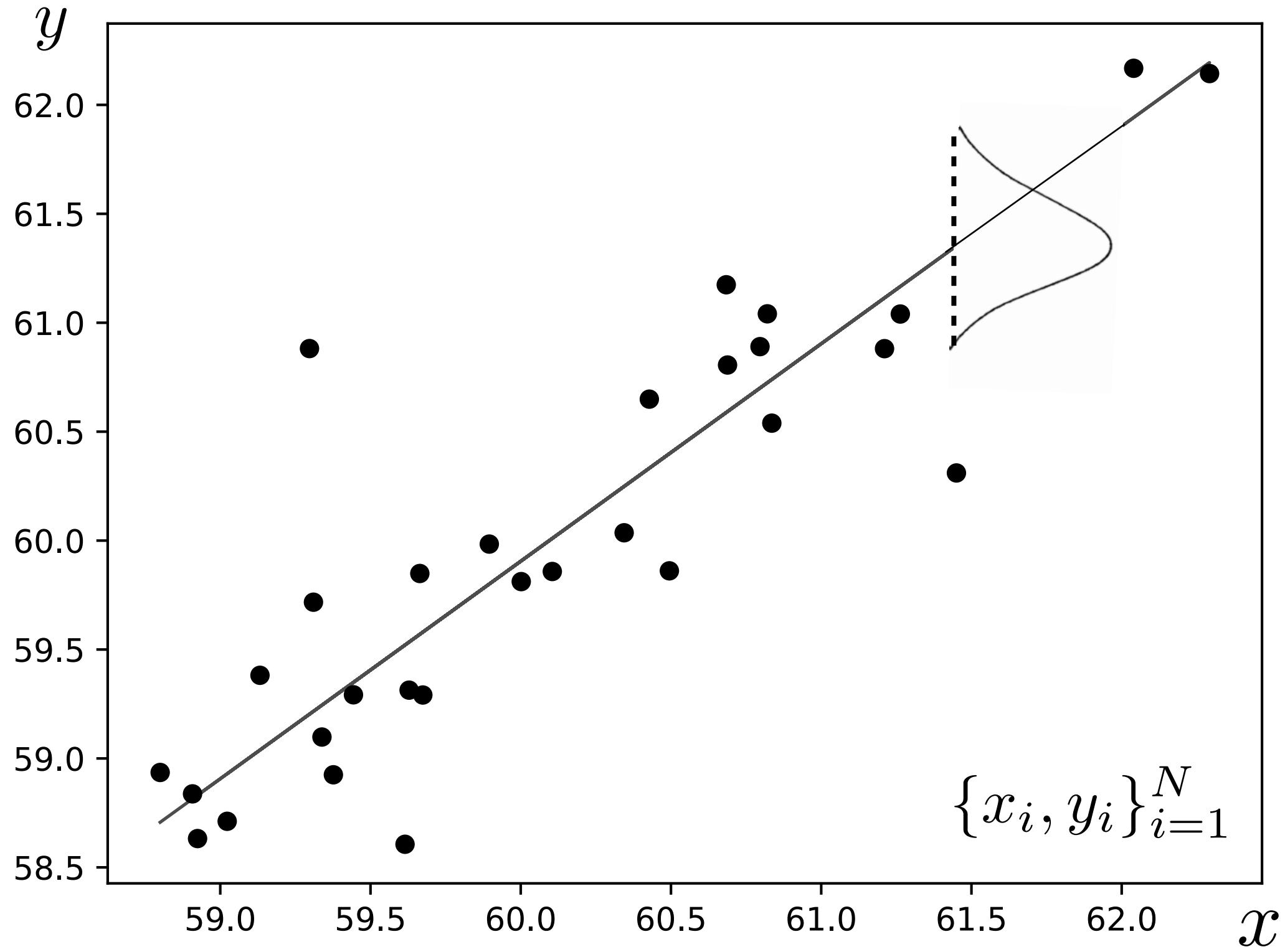
$$P_\theta(\{y_i\}_{i=1}^N | \{x_i\}_{i=1}^N) = \prod_{i=1}^N P_\theta(y_i | x_i)$$

# Why is squared error a good objective?



$$\begin{aligned}\theta^* &= \arg \max_{\theta} \prod_{I=1}^N P_{\theta}(y_i|x_i) \\ &= \arg \max_{\theta} \prod_{I=1}^N \exp -\frac{(y - f_{\theta}(x))^2}{2\sigma^2} \\ &= \arg \max_{\theta} \log \prod_{I=1}^N \exp -\frac{(y - f_{\theta}(x))^2}{2\sigma^2} \\ &= \arg \max_{\theta} \sum_{i=1}^N -\frac{(y - f_{\theta}(x))^2}{2\sigma^2} \\ &= \arg \min_{\theta} \sum_{i=1}^N \frac{(y - f_{\theta}(x))^2}{2\sigma^2} \\ &= \arg \min_{\theta} \sum_{i=1}^N (y - f_{\theta}(x))^2\end{aligned}$$

# Why is squared error a good objective?



Under the assumption that the data is distributed as:

$$Y = f_\theta(X) + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

$$P_\theta(\{y_i\}_{i=1}^N | \{x_i\}_{i=1}^N) = \prod_{i=1}^N P_\theta(y_i | x_i)$$

The most probable estimate of  $\theta$  is:

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

This is called the **maximum likelihood** estimator of  $\theta$ .

# How to minimize the objective w.r.t. $\theta$ ?

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

We'll use methods from **optimization**.

# How to minimize the objective w.r.t. $\theta$ ?

In the linear case:

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

Learning problem

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

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

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

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

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

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

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

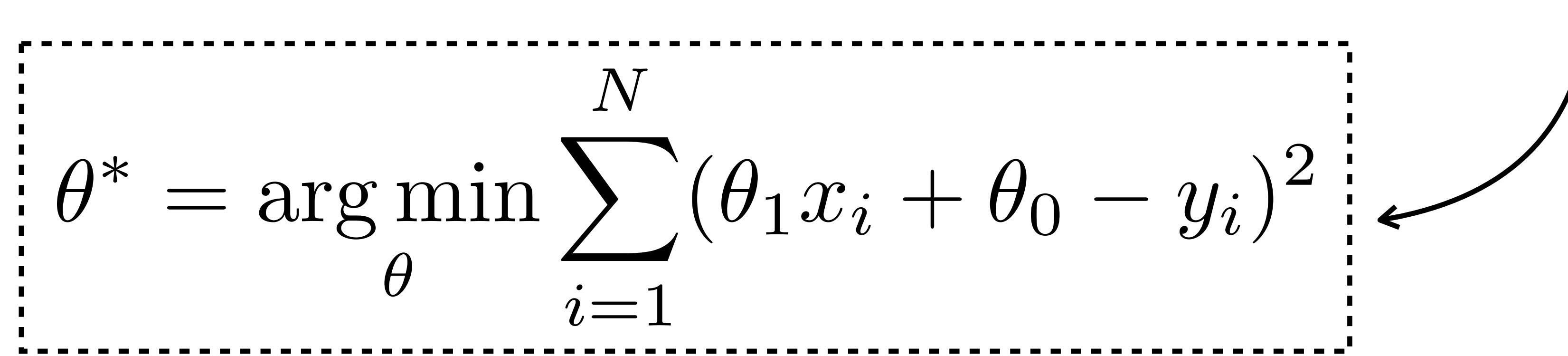
$$\boxed{\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

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)

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

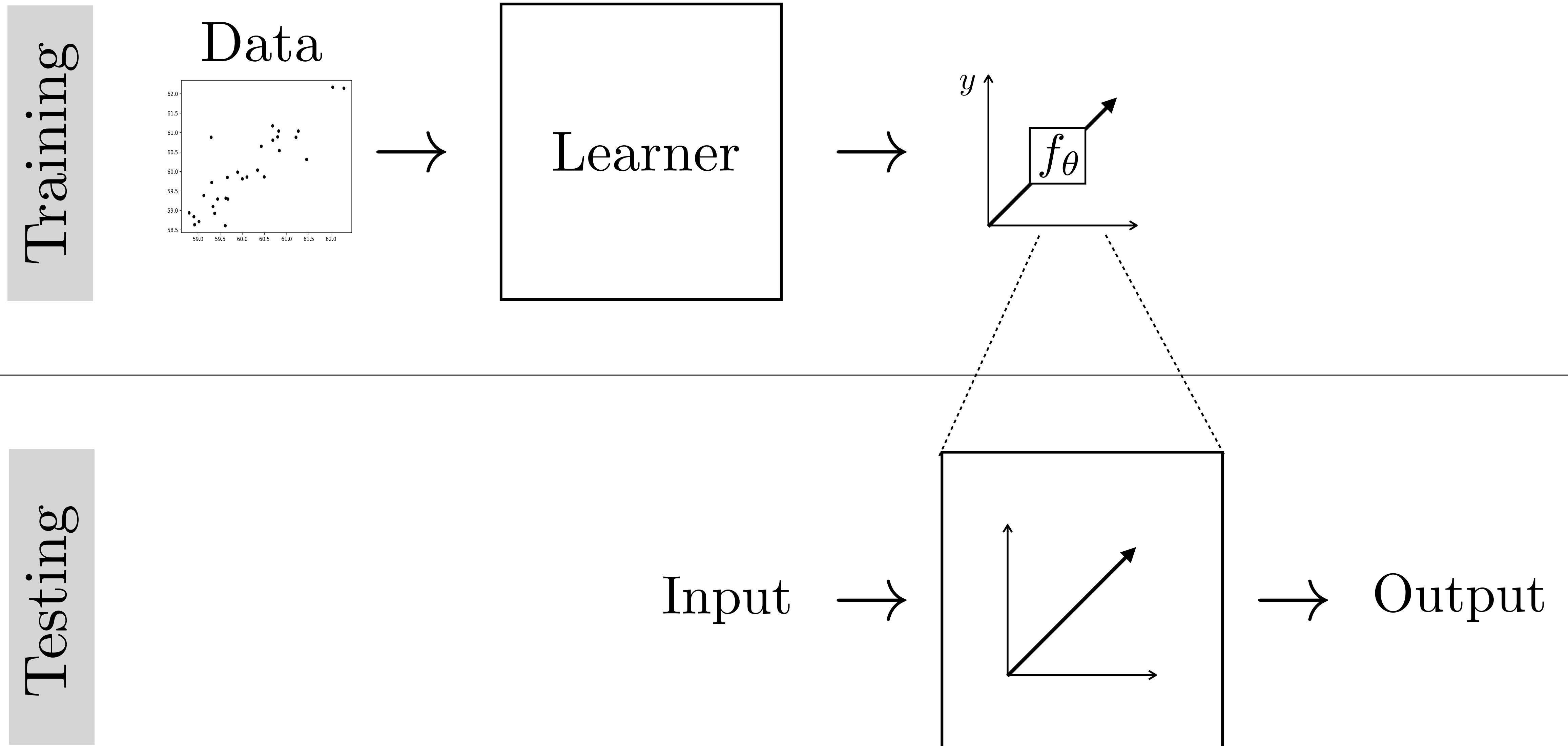
Objective function  
(loss)

Hypothesis space

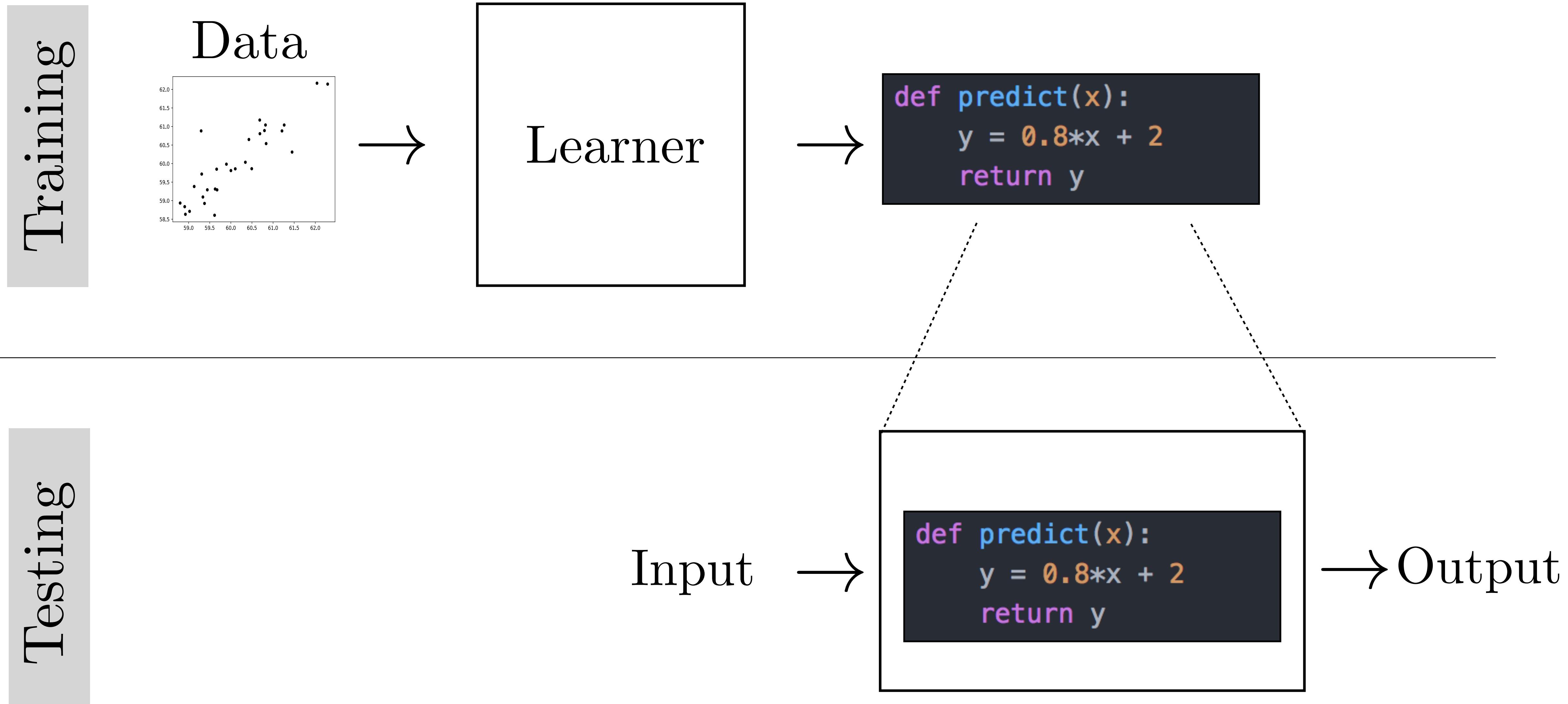
Training data

The diagram illustrates the components of the Empirical Risk Minimization formula. It shows the formula  $f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i), \mathbf{y}_i)$ . Three annotations are provided: 'Hypothesis space' with an arrow pointing to the  $f \in \mathcal{F}$  term; 'Training data' with two arrows pointing to the summation symbol ( $\sum$ ) and the arguments of the loss function ( $f(\mathbf{x}_i)$  and  $\mathbf{y}_i$ ); and 'Objective function (loss)' with an arrow pointing to the  $\mathcal{L}$  term.

# Example 1: Linear least squares



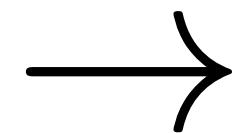
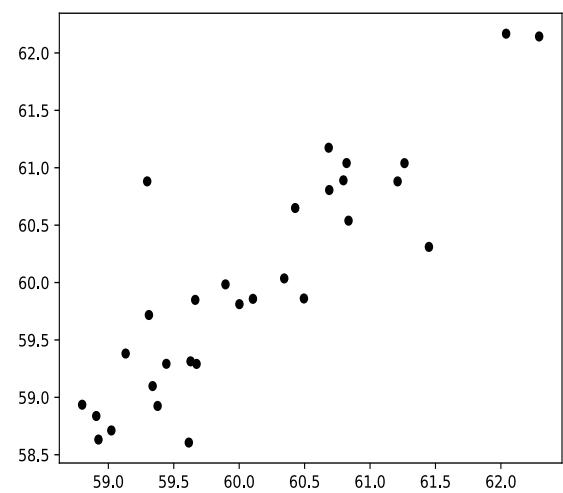
# Example 2: Program induction



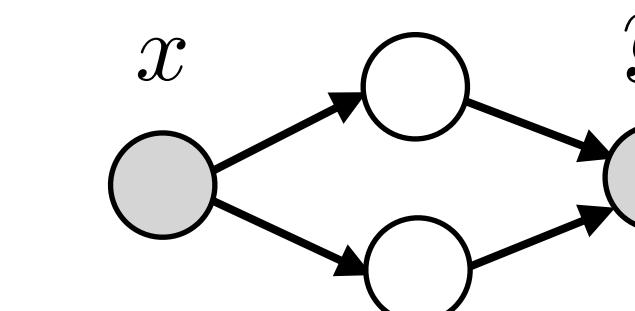
# Example 3: Deep learning

Training

Data

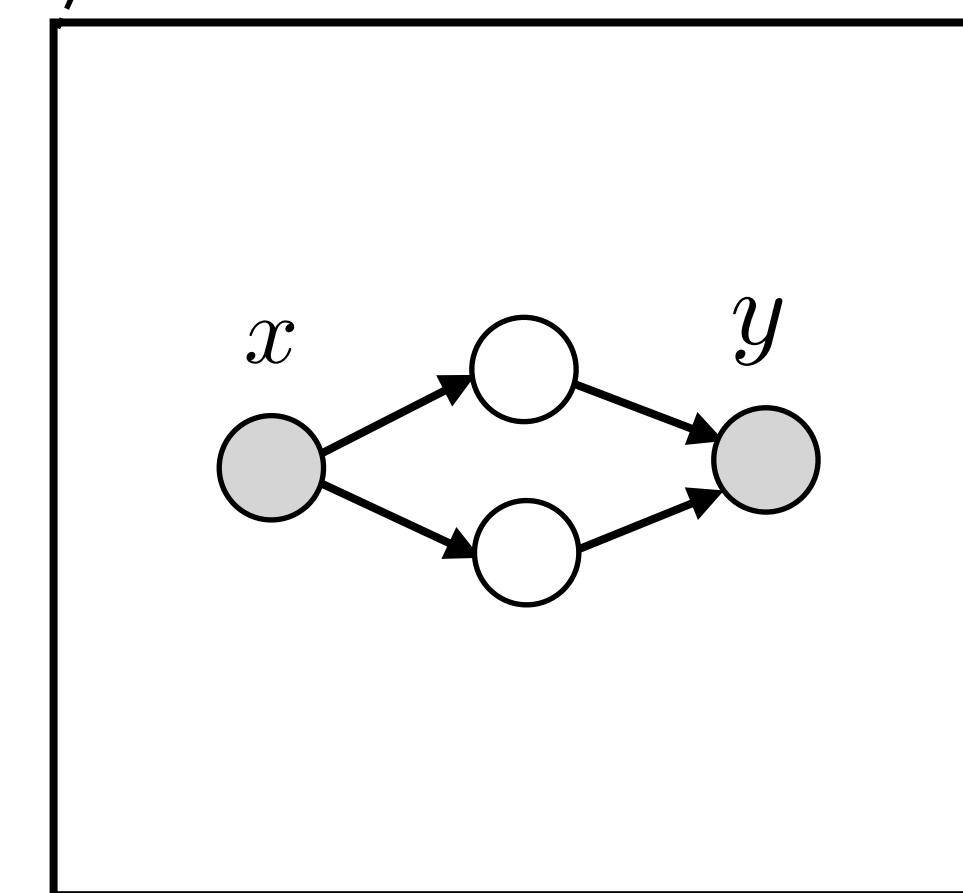


Learner



Testing

Input →



→ Output

# Learning for vision

Questions:

1. How do you represent the input and output?
2. What is the objective?
3. What is the hypothesis space? (e.g., linear, polynomial, neural net?)
4. How do you optimize? (e.g., gradient descent, Newton's method?)
5. What data do you train on?

# Learning for vision

Questions:

- 1. How do you represent the input and output?**
- 2. What is the objective?**
3. What is the hypothesis space? (e.g., linear, polynomial, neural net?)
4. How do you optimize? (e.g., gradient descent, Newton's method?)
5. What data do you train on?

# Image classification

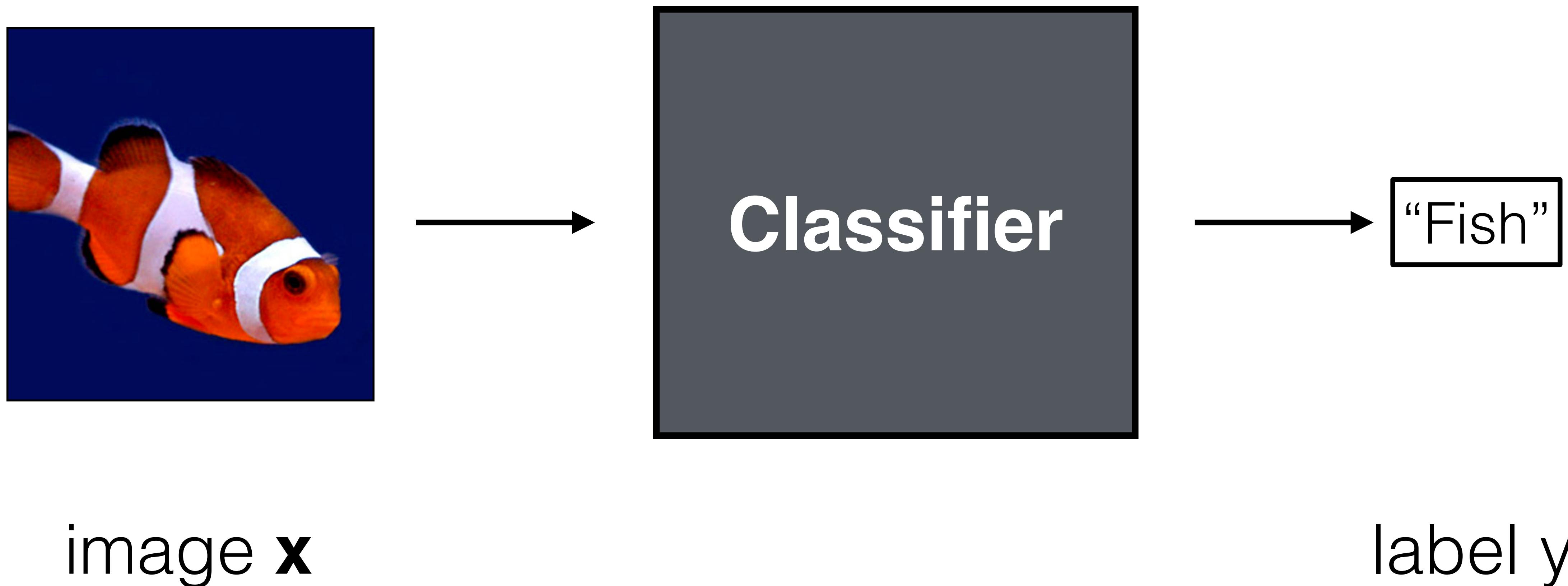


image x

label y

# Image classification



“Fish”

image  $x$

label  $y$

# Image classification

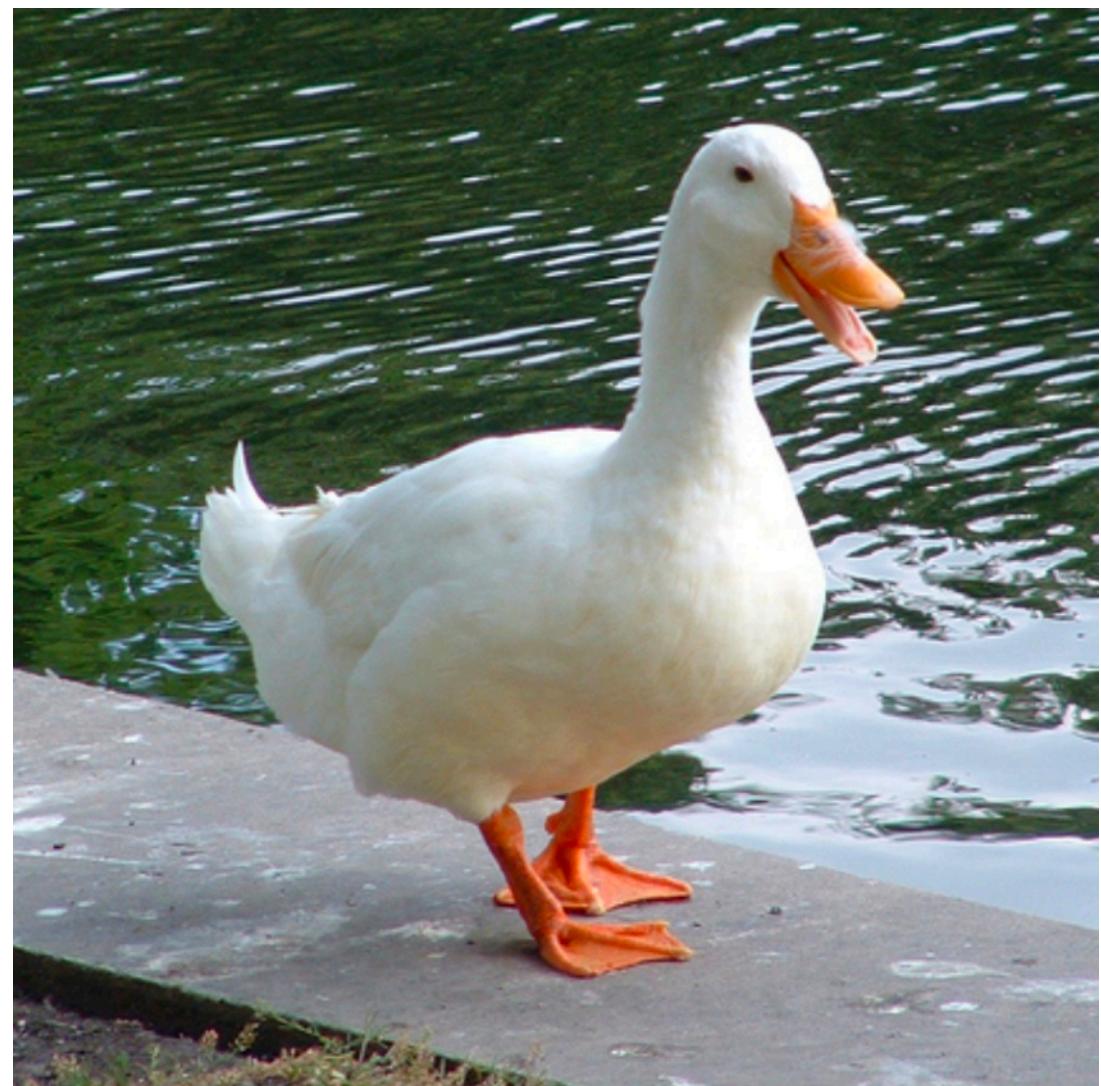


“Fish”

image  $x$

label  $y$

# Image classification



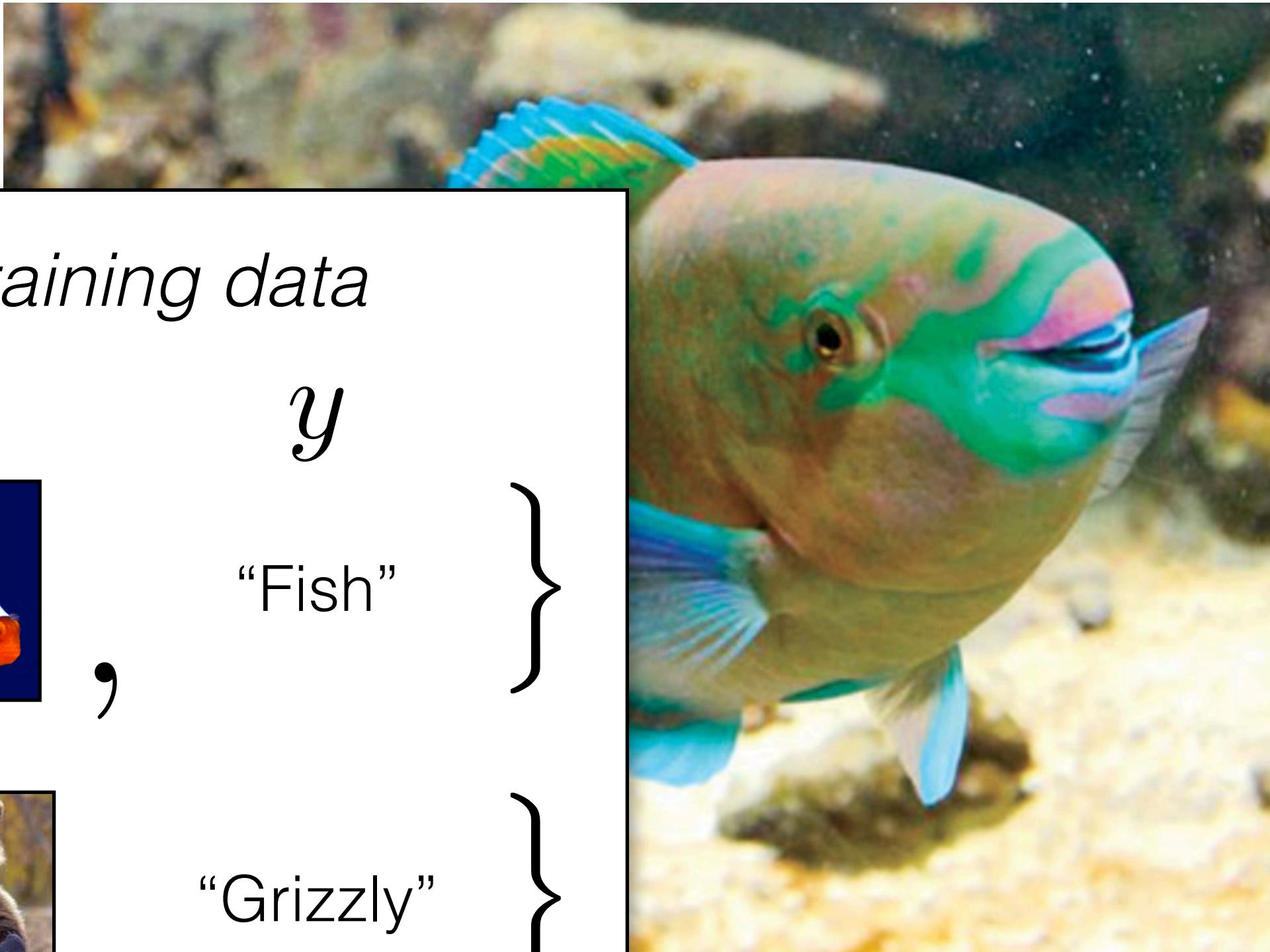
“Duck”

:

image **x**

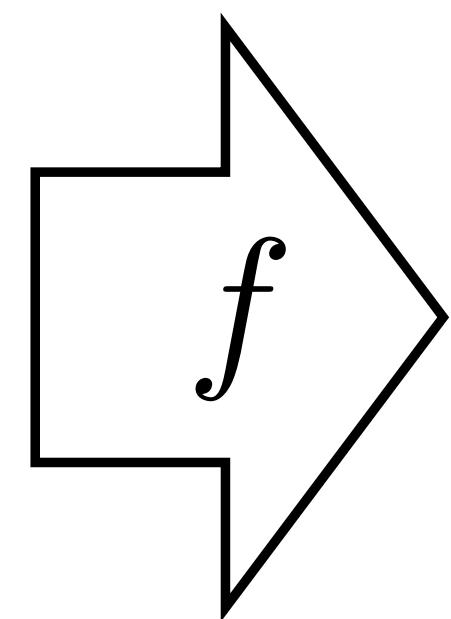
label **y**

**X**



*Training data*

<b>x</b>	<b>y</b>
{  , }	“Fish”
{  , }	“Grizzly”
{  , }	“Chameleon”
:	

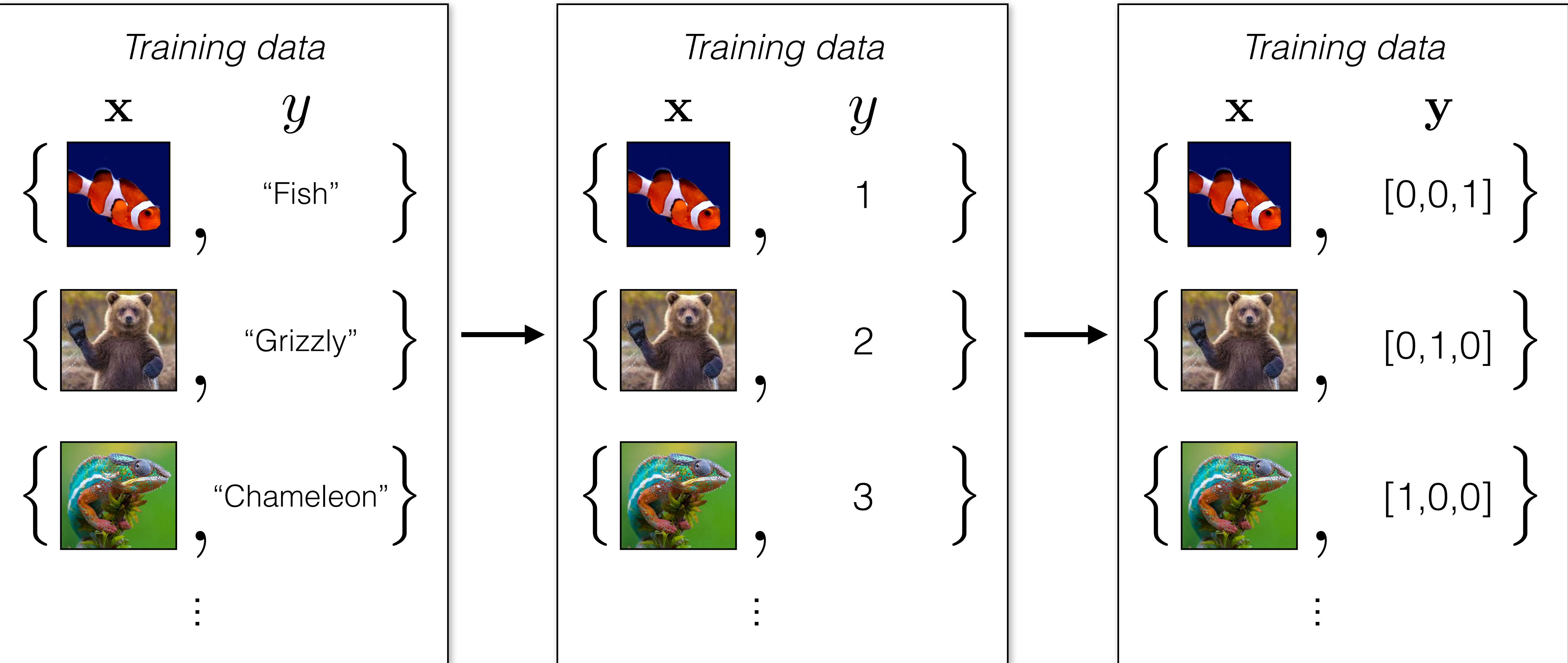


**y**  
“Fish”

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

# How to represent class labels?

**One-hot vector**



# What should the loss be?

**0-1 loss** (number of misclassifications)

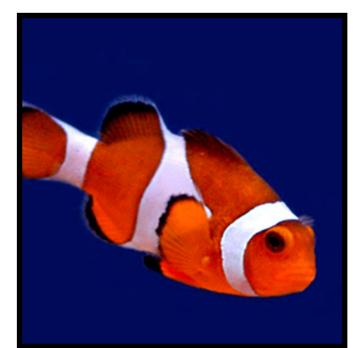
$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \mathbb{1}(\hat{\mathbf{y}} = \mathbf{y}) \quad \leftarrow \text{discrete; hard to optimize!}$$

**Cross entropy**

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad \leftarrow \begin{array}{l} \text{continuous,} \\ \text{differentiable,} \\ \text{convex} \end{array}$$

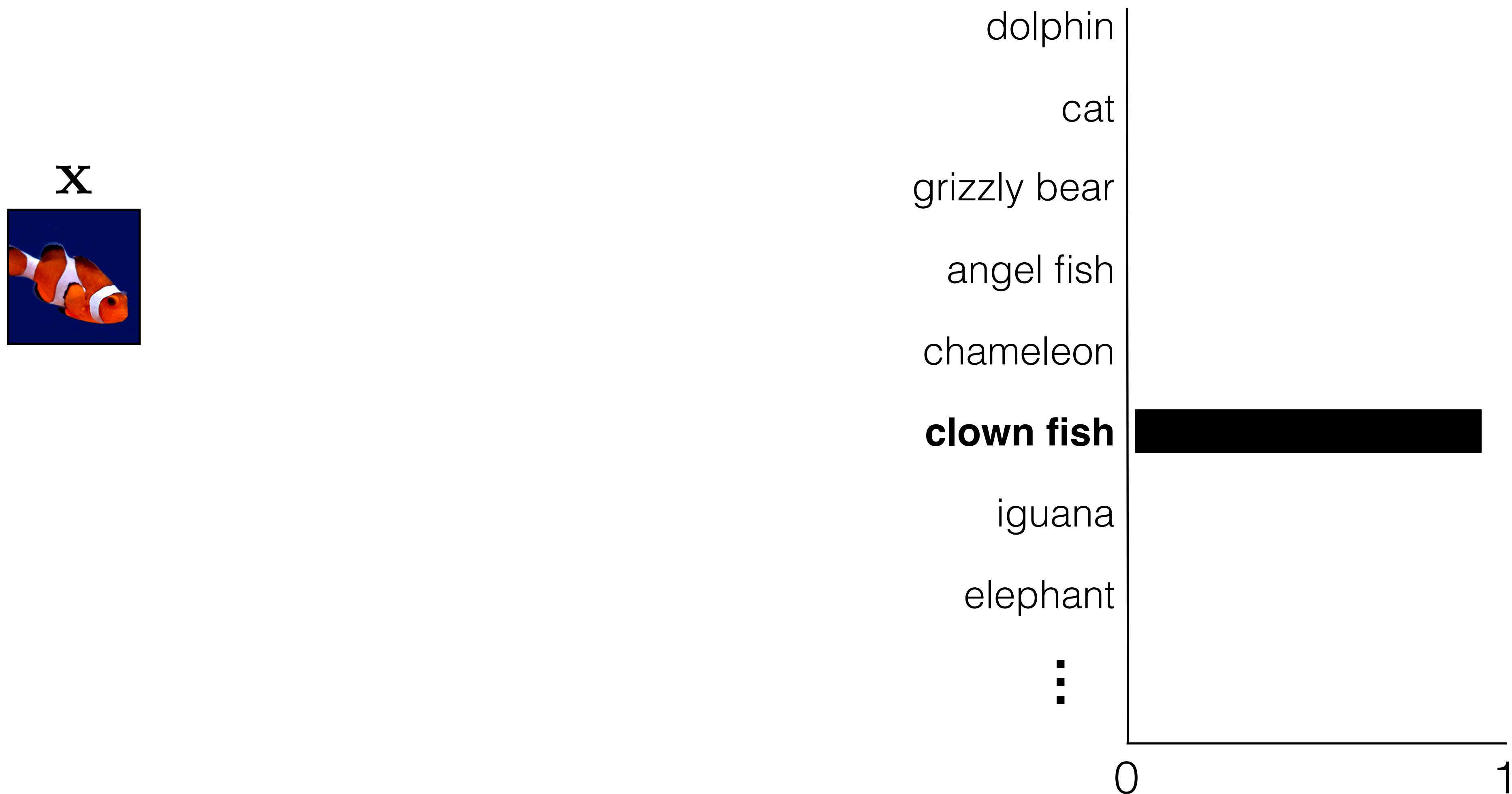
Ground truth label  $y$

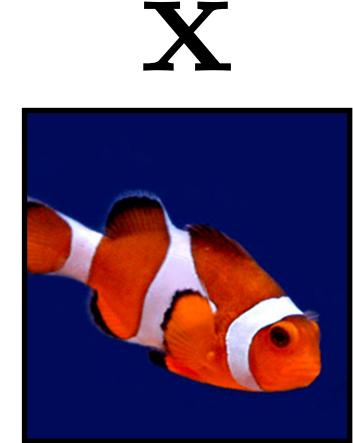
$x$



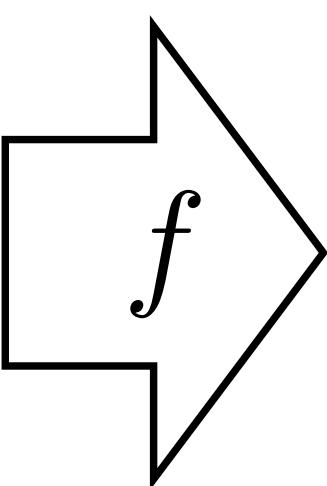
[0,0,0,0,0,1,0,0,...]

Ground truth label **y**





**X**

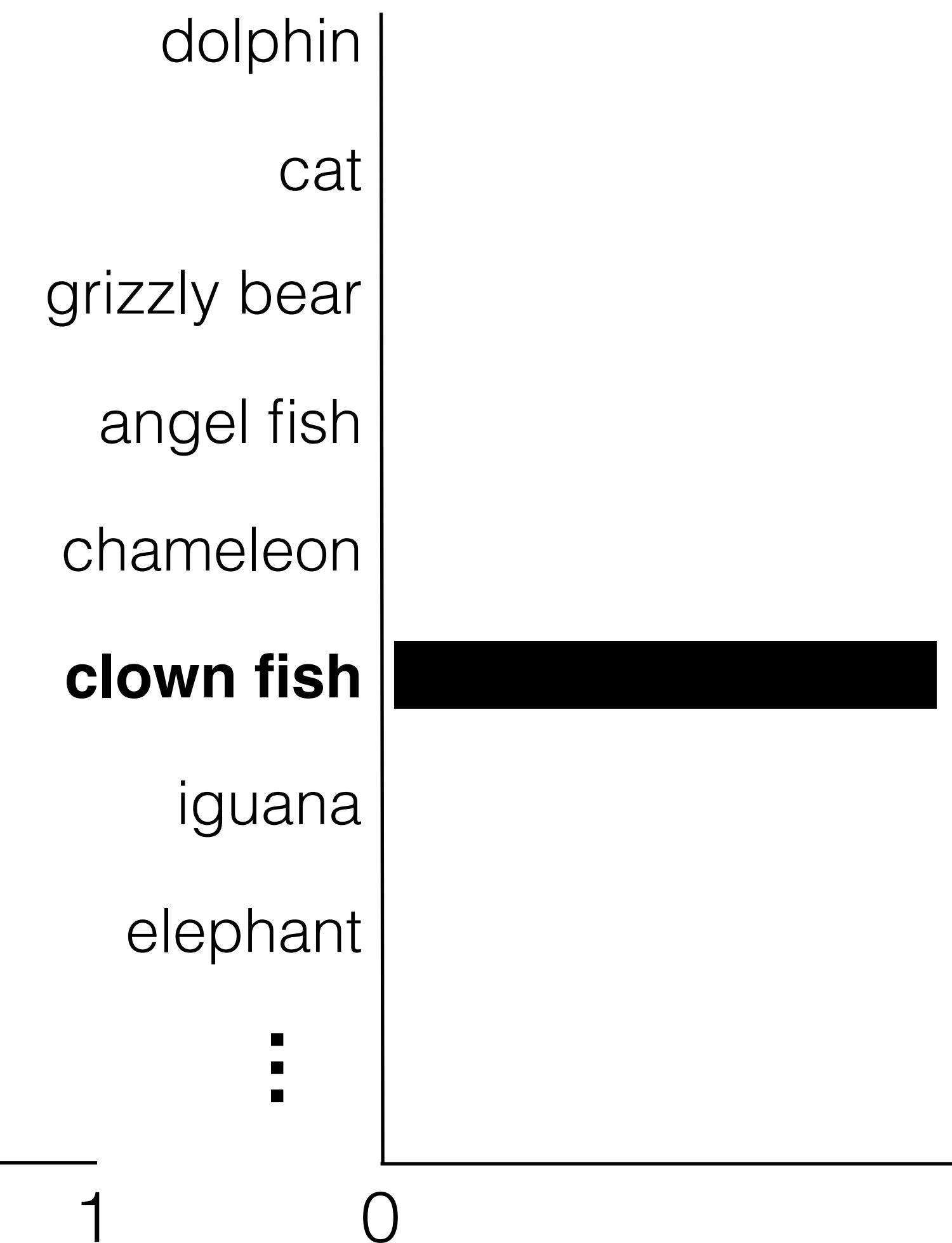


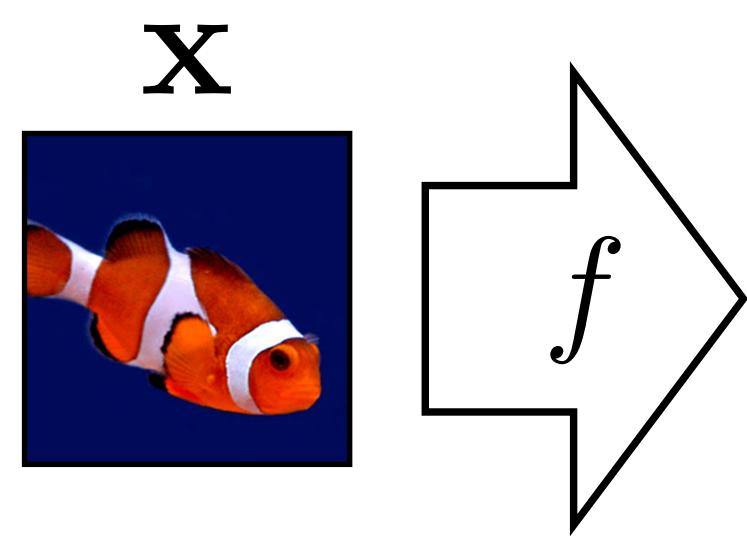
Prediction  $\hat{\mathbf{y}}$

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



Ground truth label  $\mathbf{y}$

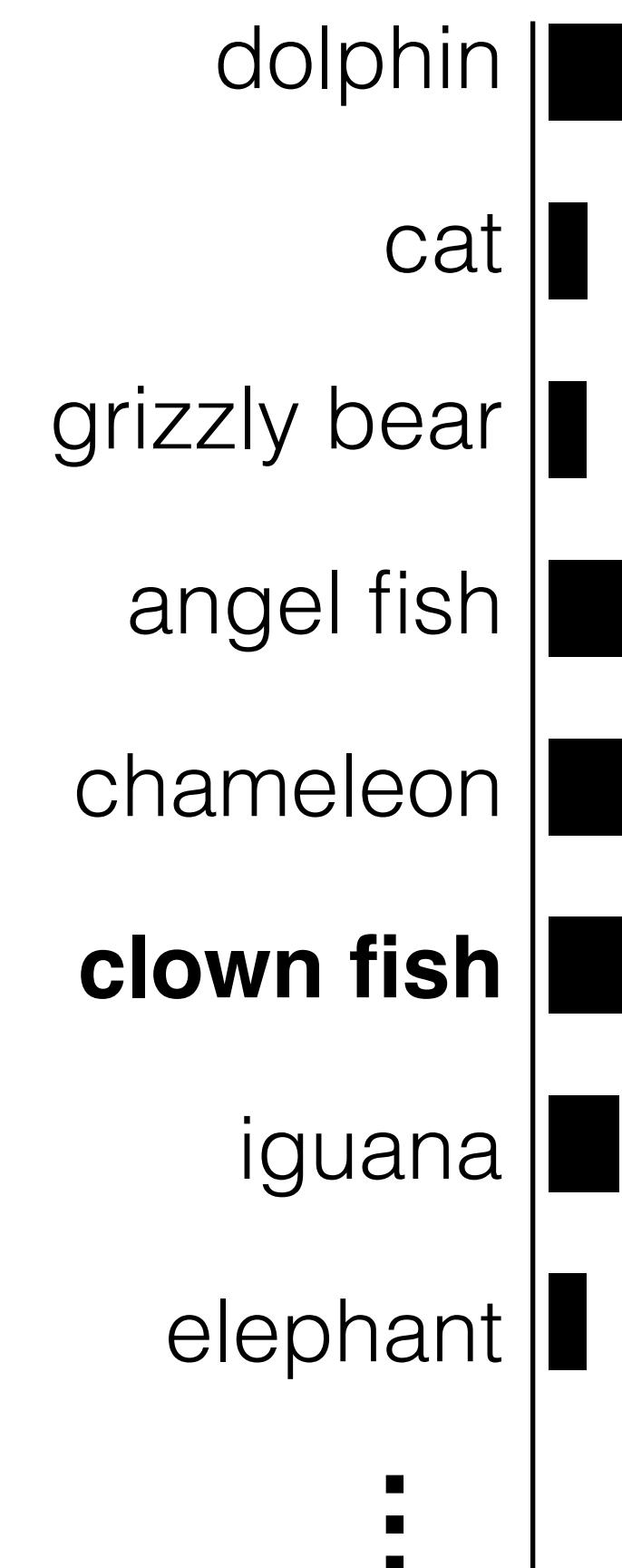




**X**

Prediction  $\hat{\mathbf{y}}$

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

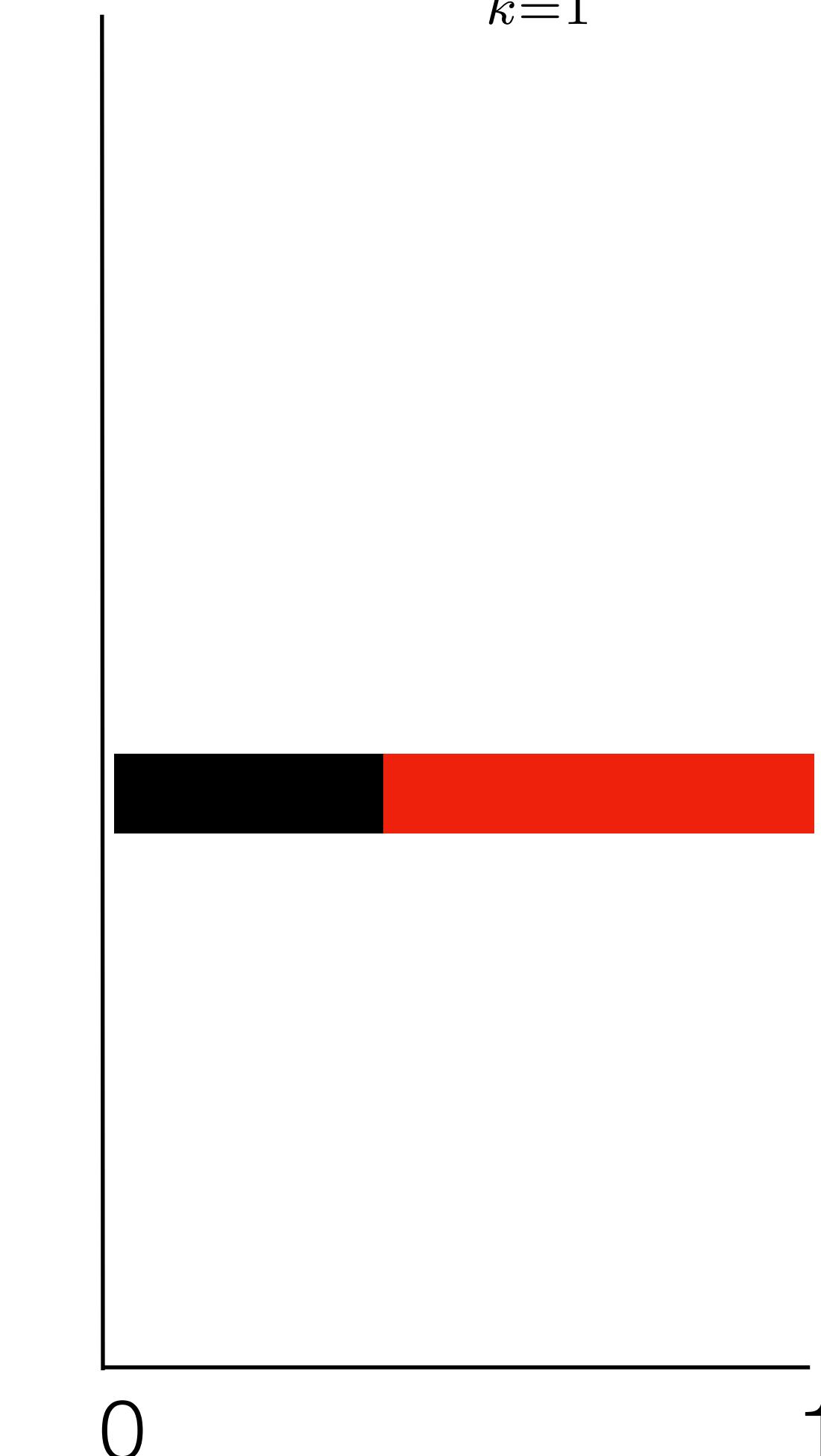


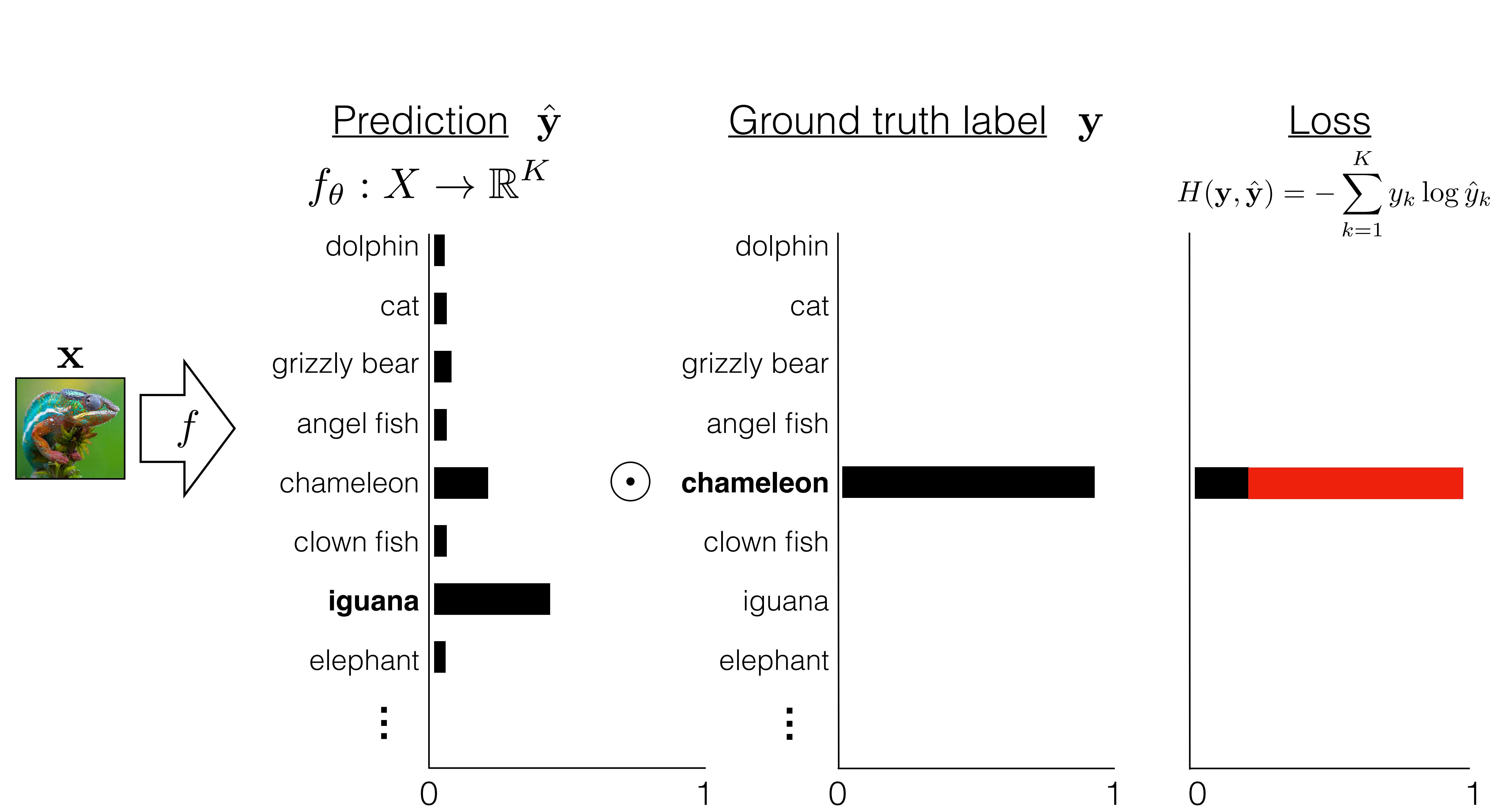
Ground truth label  $\mathbf{y}$



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)

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

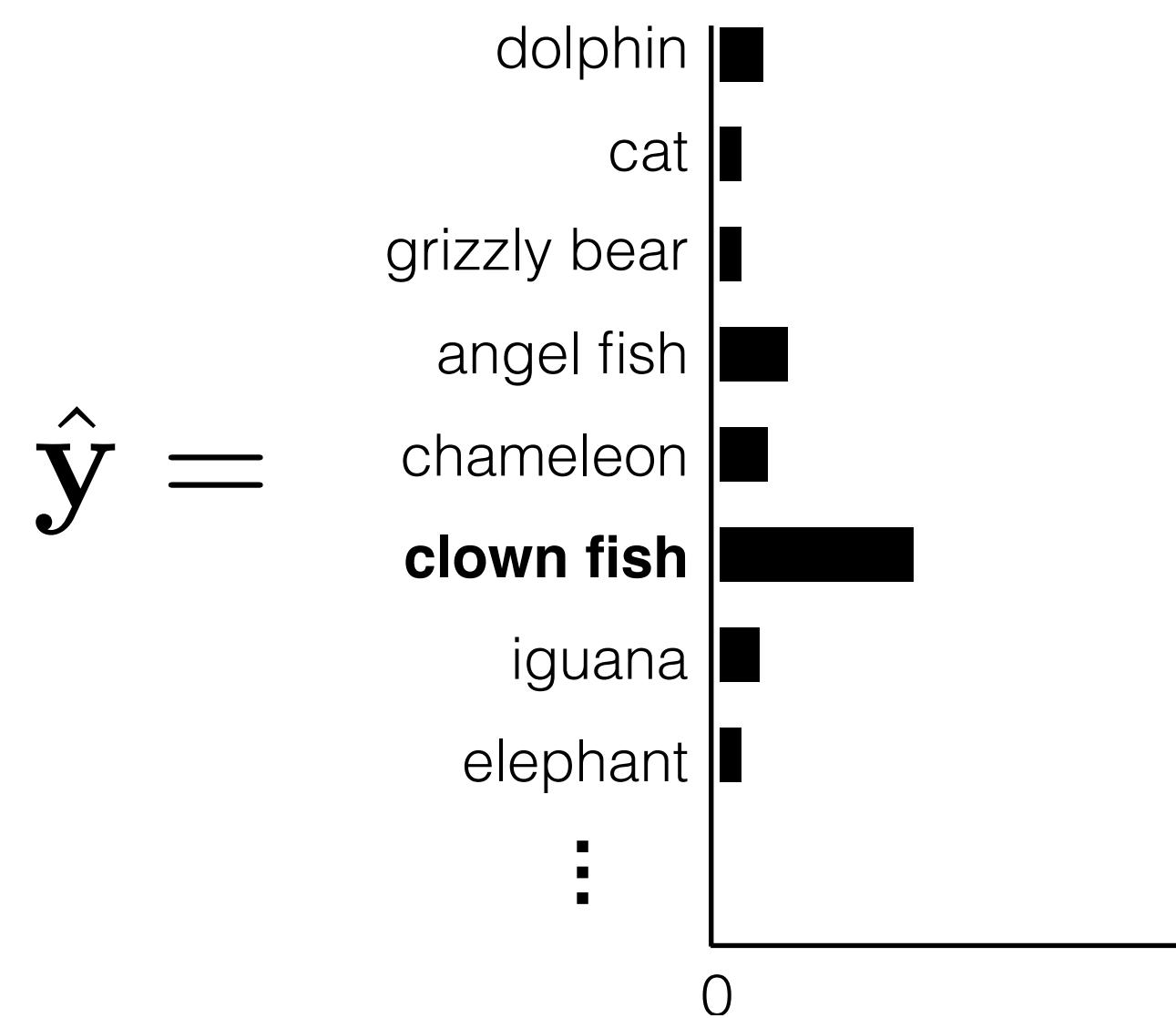
$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

← **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 mass function!**

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



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

Probabilistic interpretation:

$$\hat{\mathbf{y}} \equiv [P_{\theta}(Y = 1|X = \mathbf{x}), \dots, P_{\theta}(Y = K|X = \mathbf{x})] \leftarrow \text{predicted probability of each class given input } \mathbf{x}$$

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \leftarrow \text{picks out the -log likelihood of the ground truth class } \mathbf{y \text{ under the model prediction } \hat{\mathbf{y}}}$$

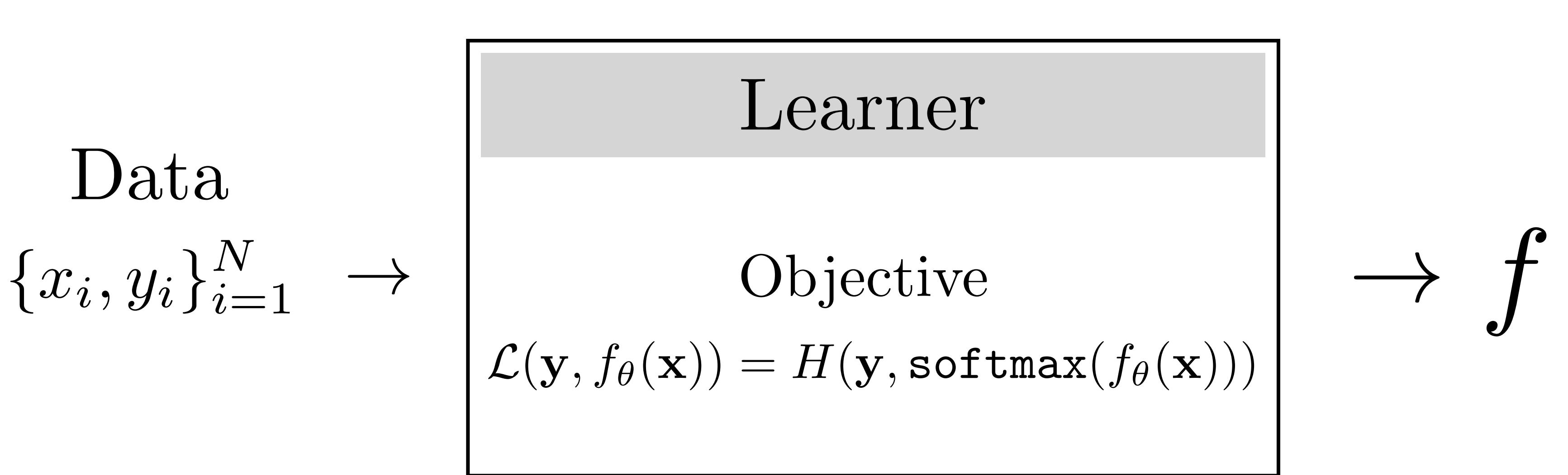
$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N H(\mathbf{y}_i, \hat{\mathbf{y}}_i) \leftarrow \text{maximum likelihood}$$

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

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

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

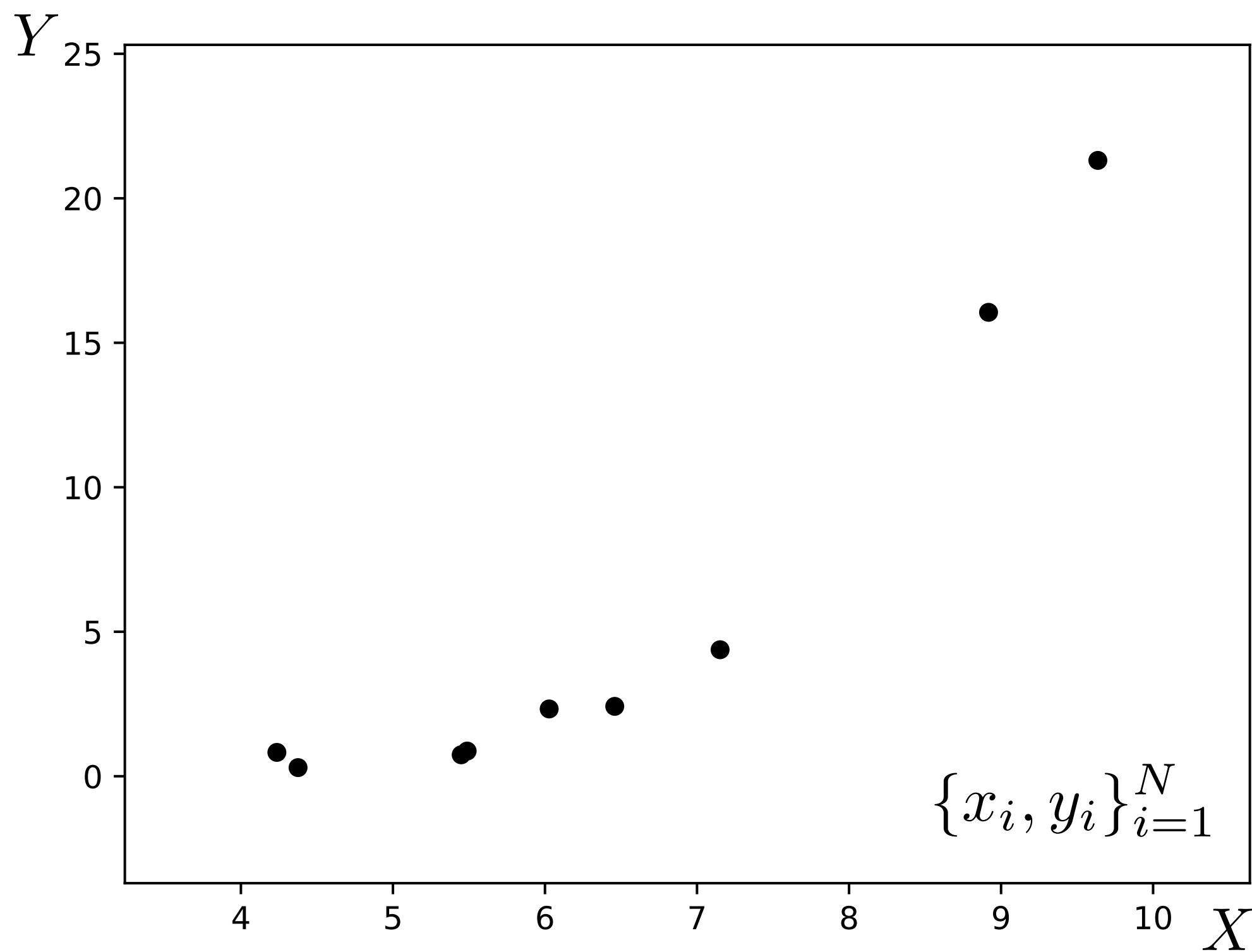
$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$



# The Problem of Generalization

# Linear regression

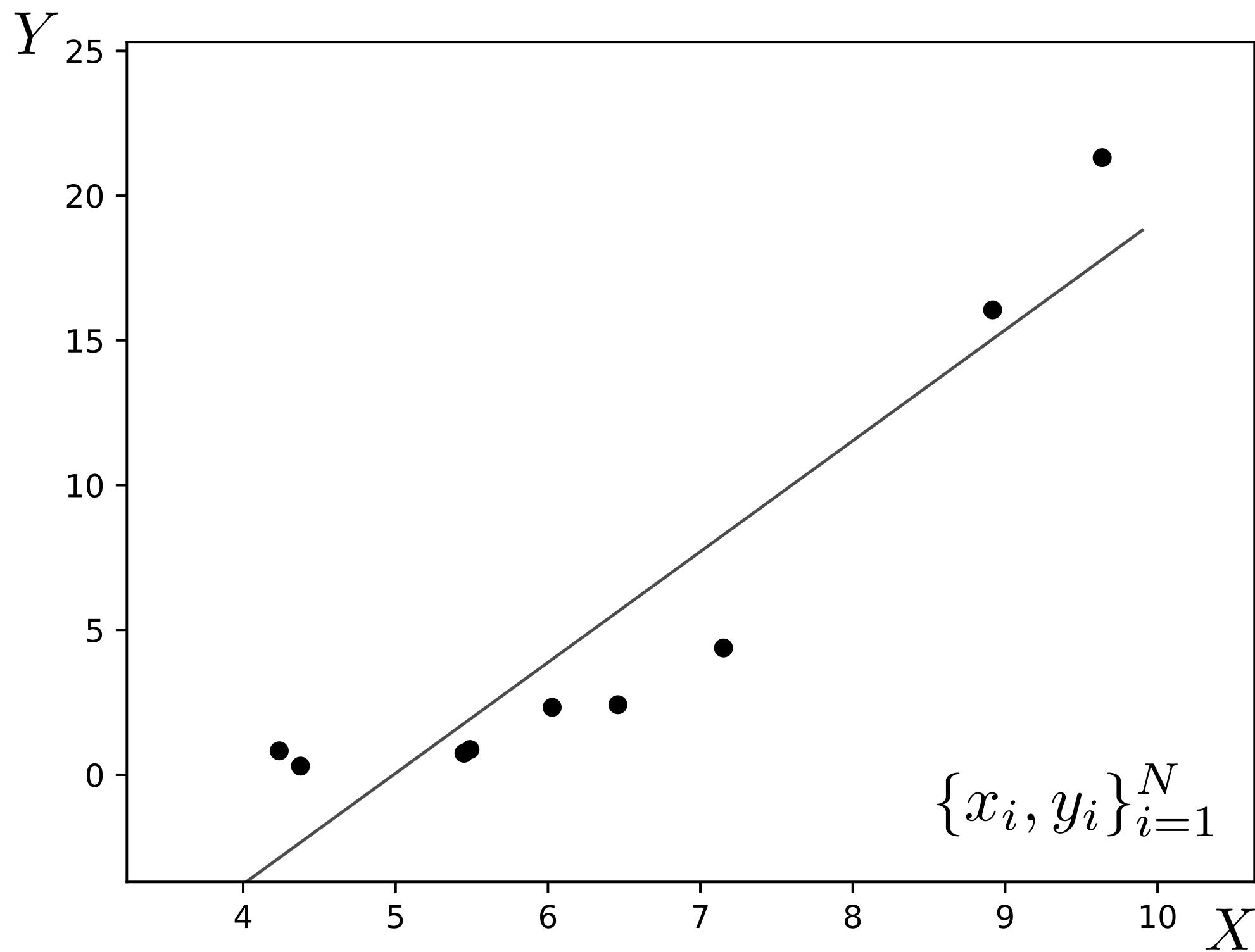
Training data



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

# Linear regression

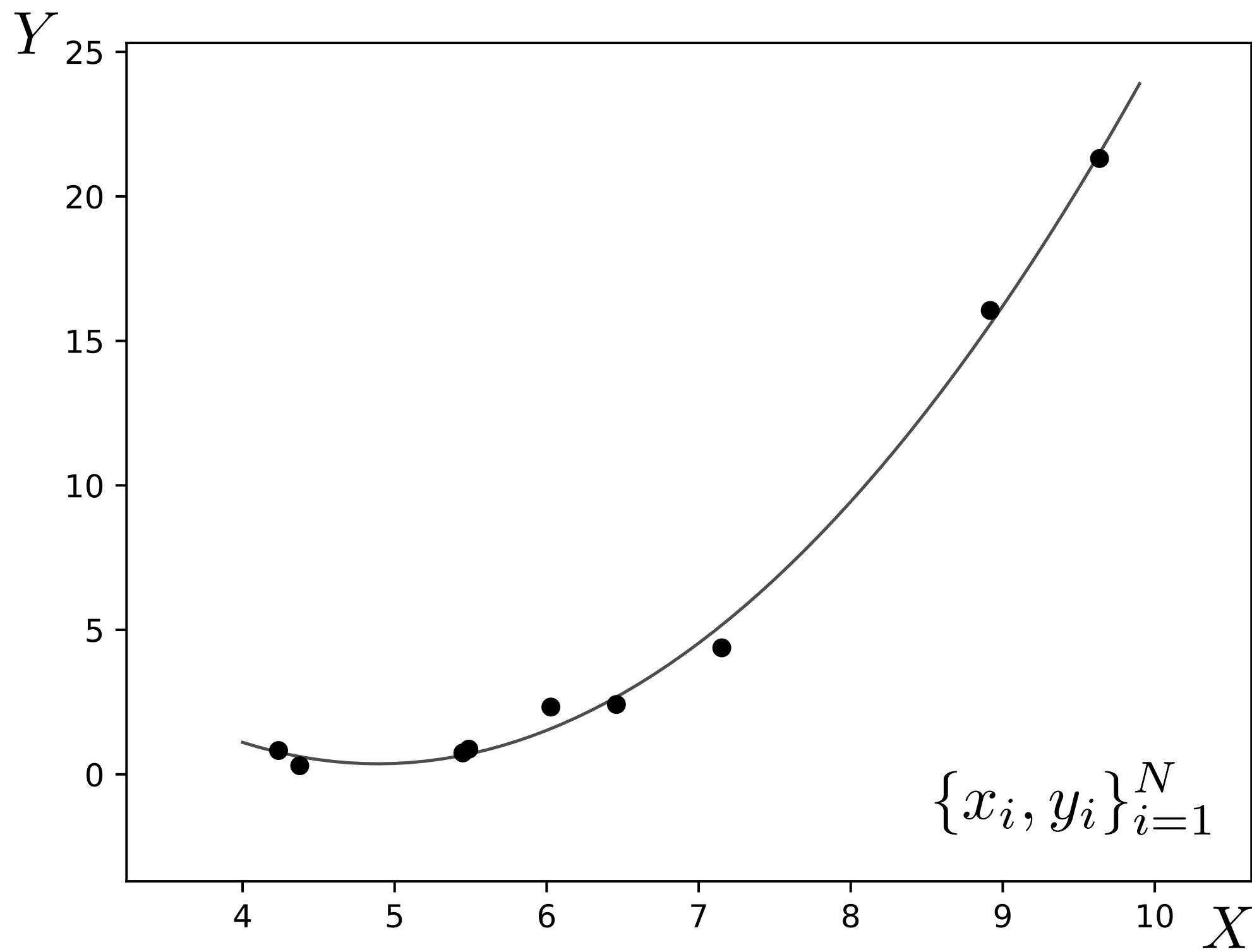
Training data



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

# Polynomial regression

Training data

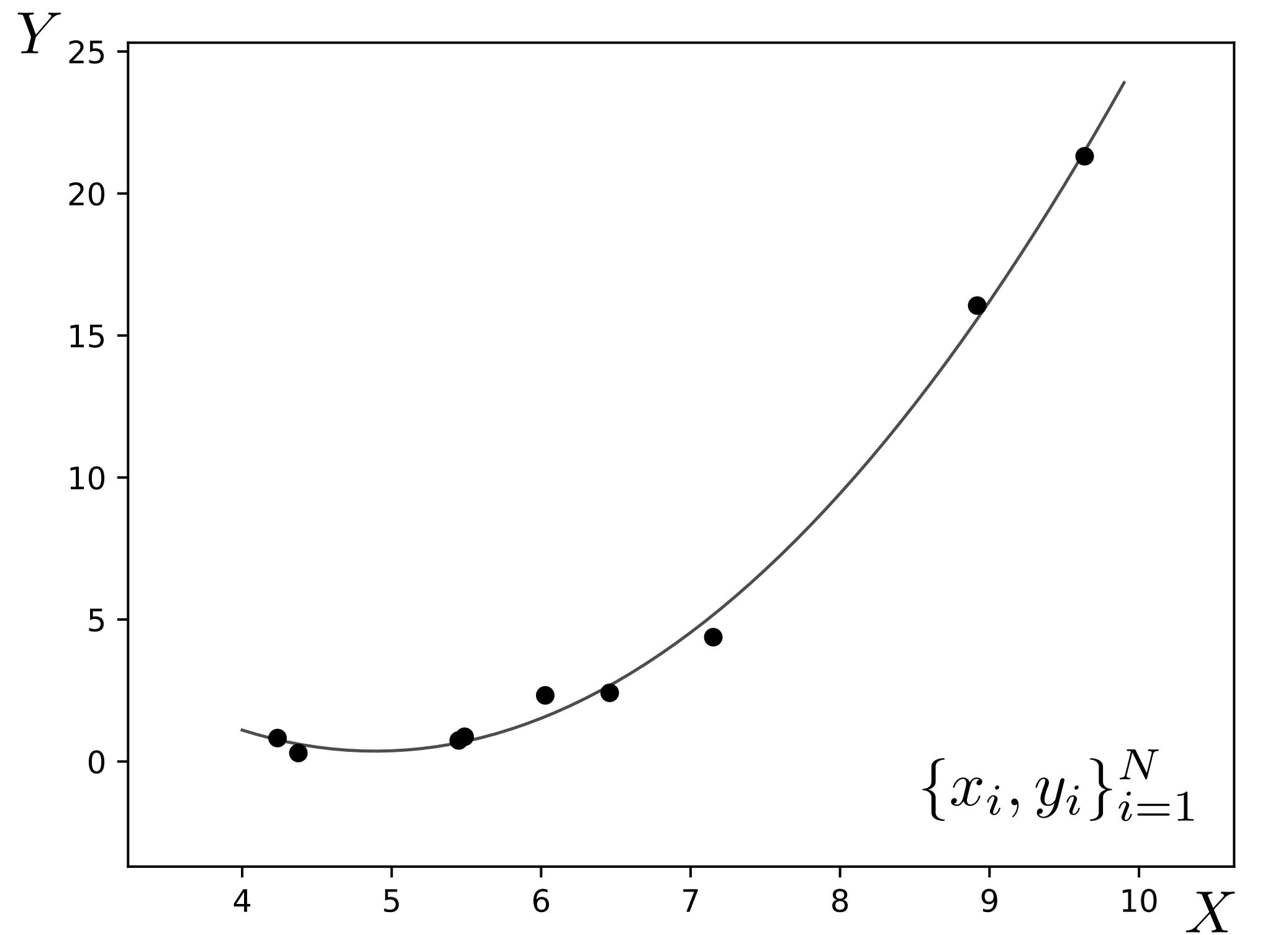


$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

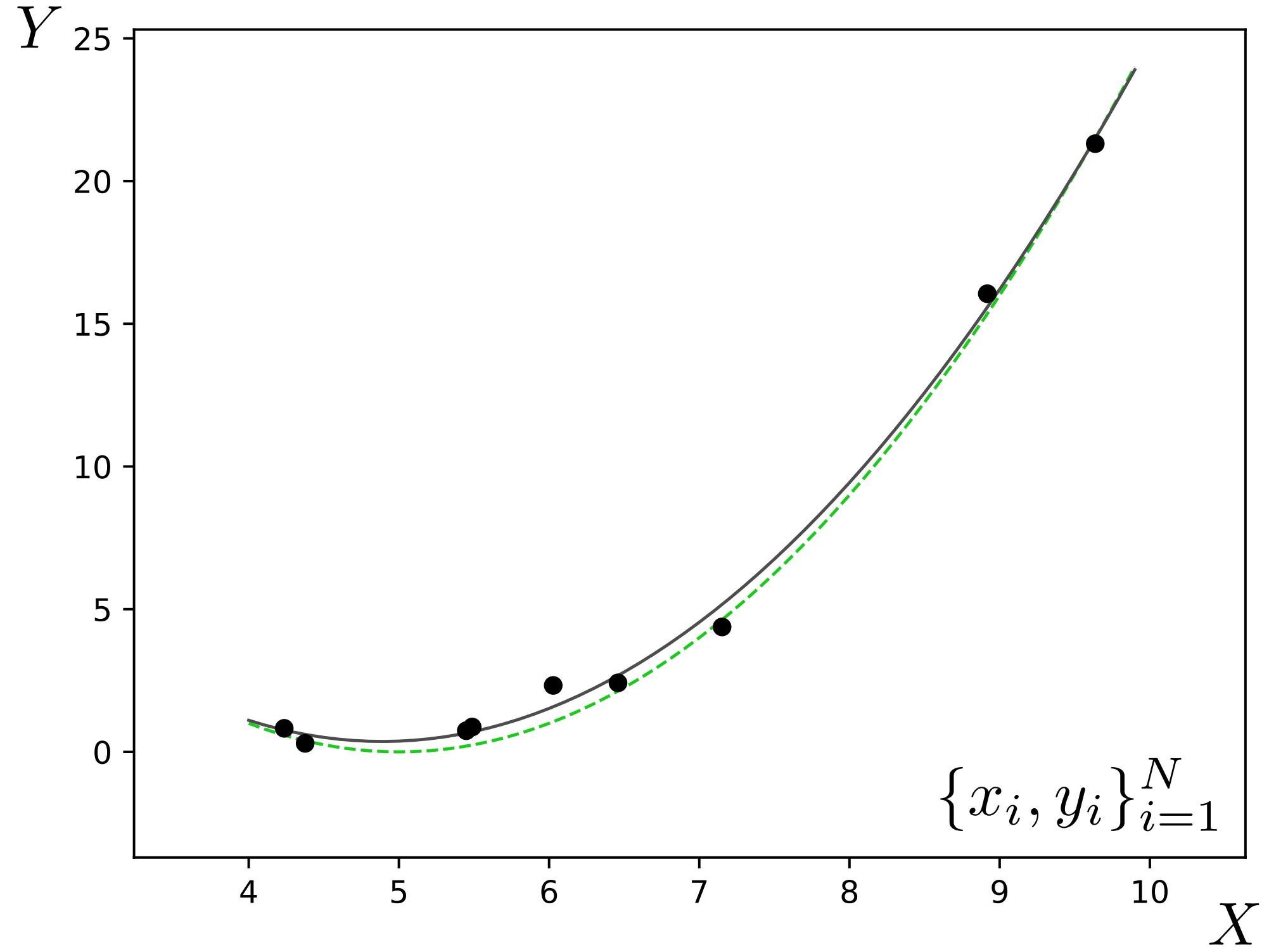
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

K-th degree polynomial regression

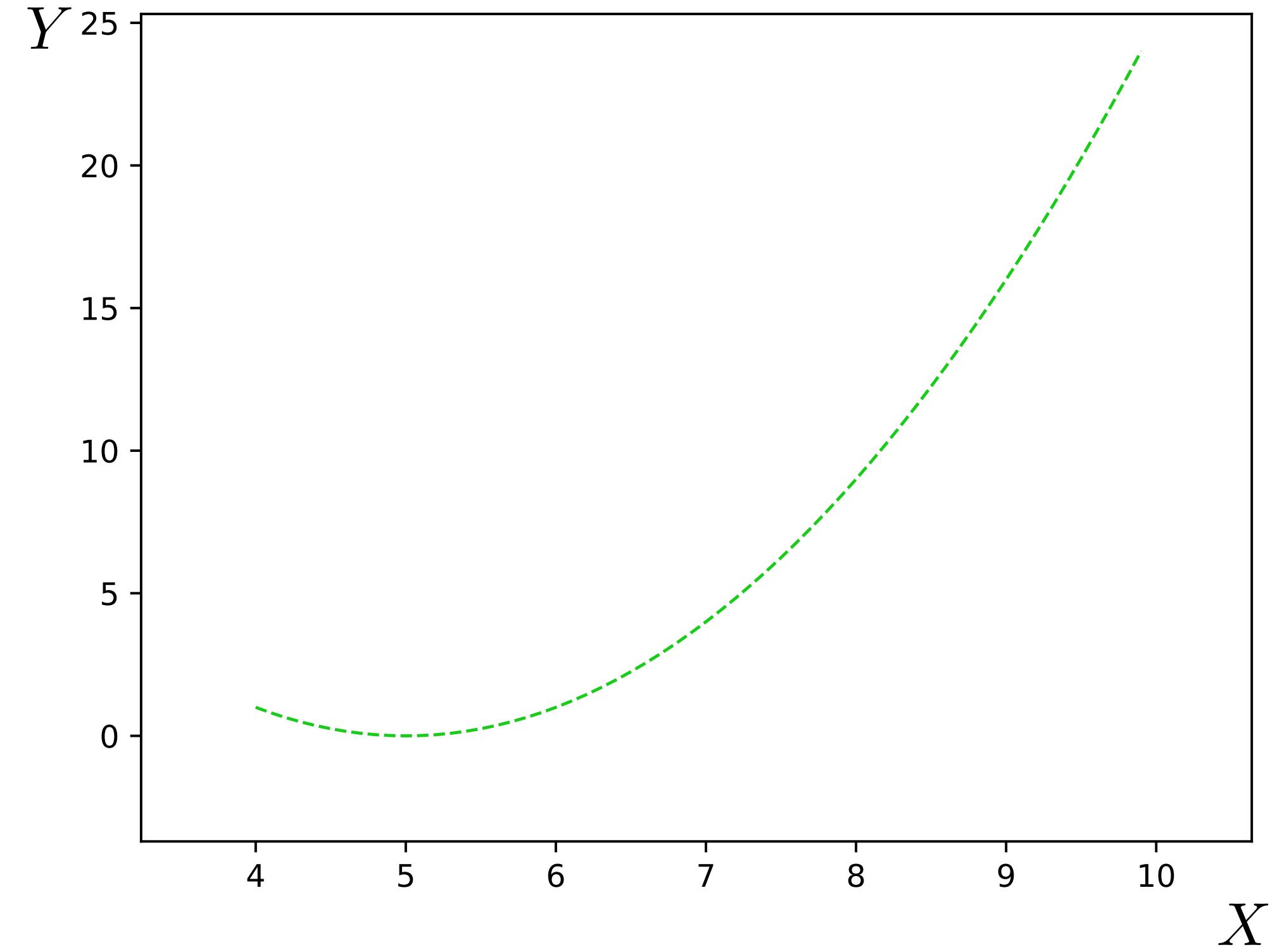
# Training data



## Training data



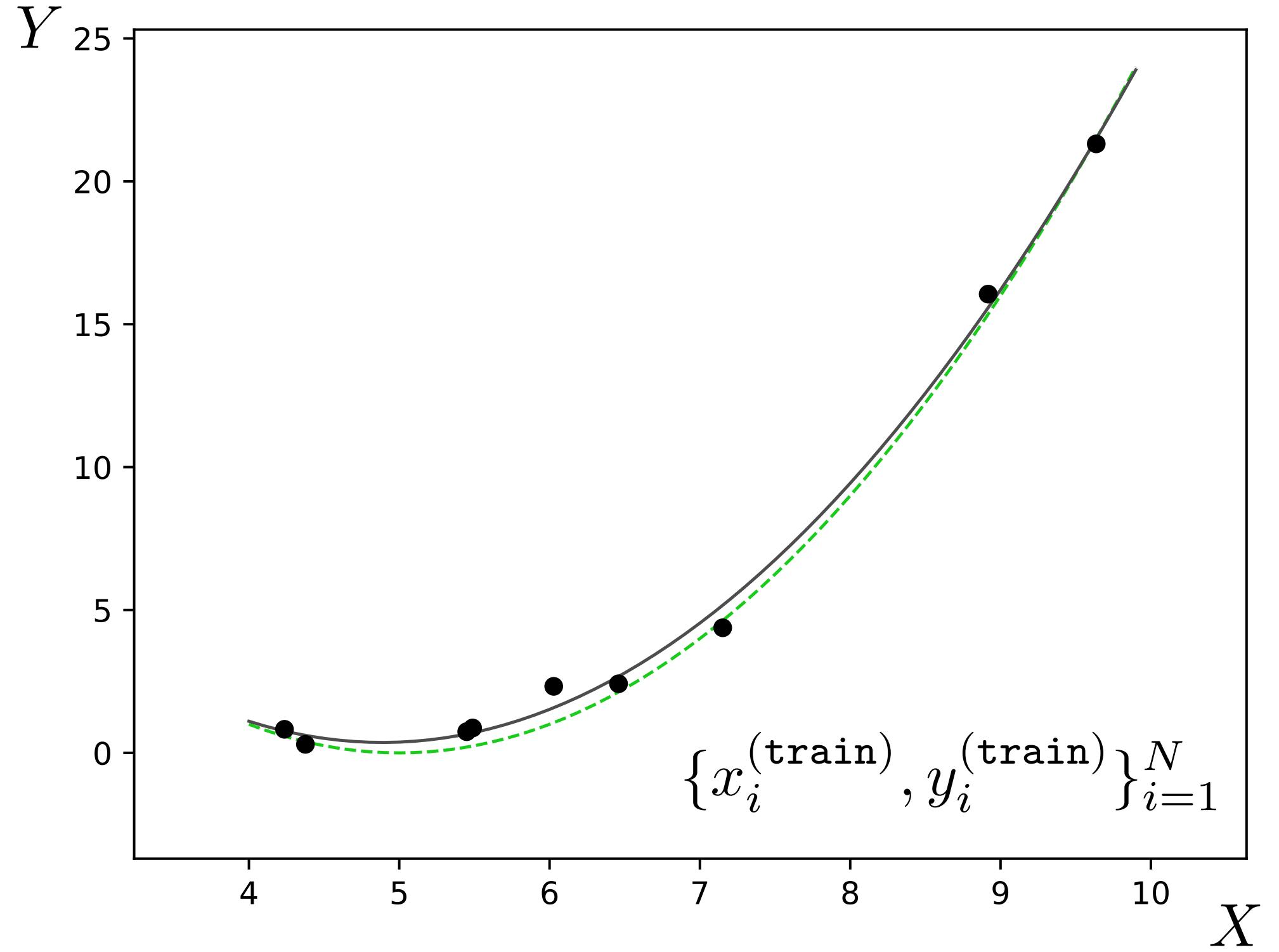
## Test data



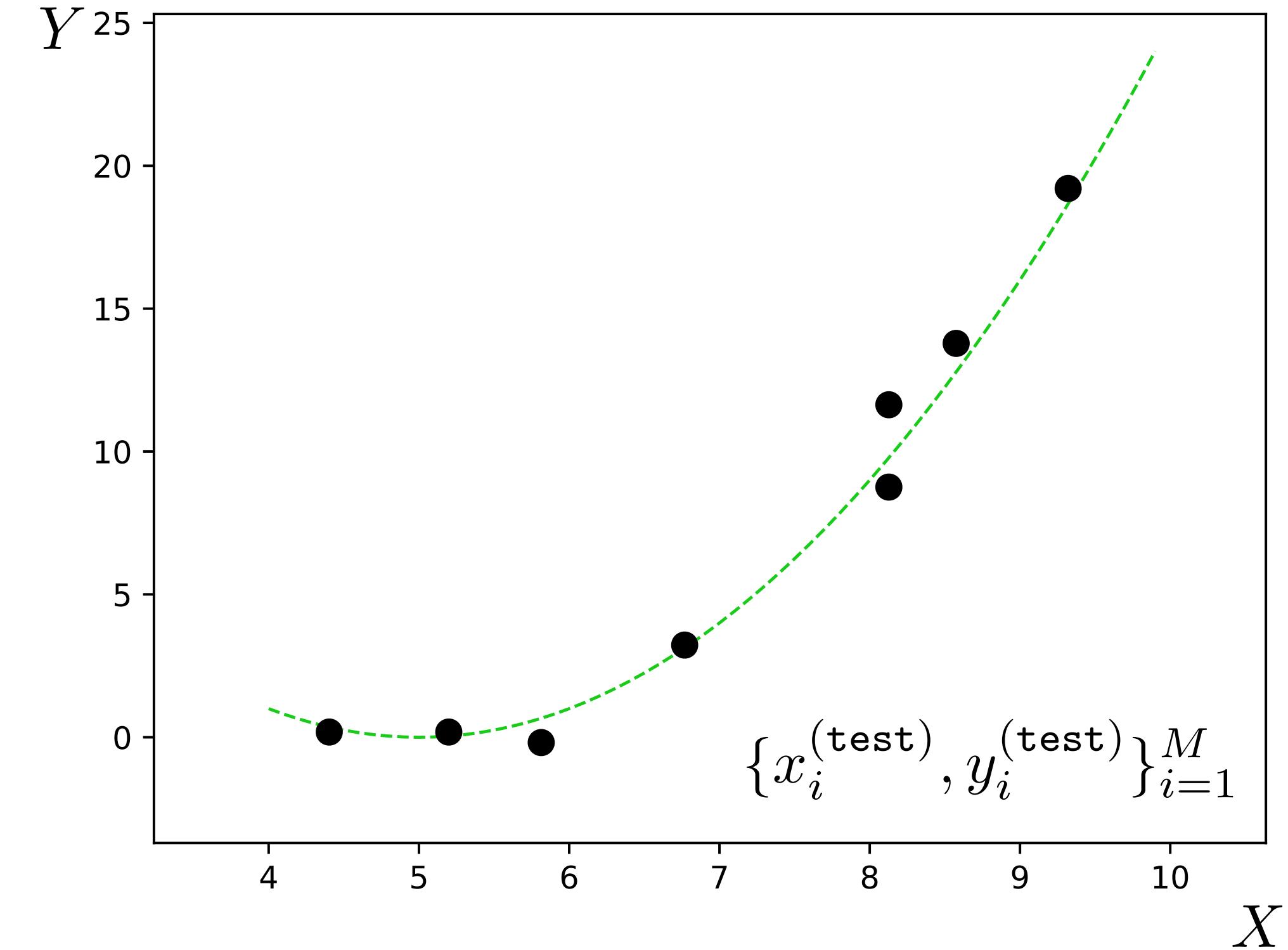
True data-generating process

$p_{\text{data}}$

## Training data



## Test data



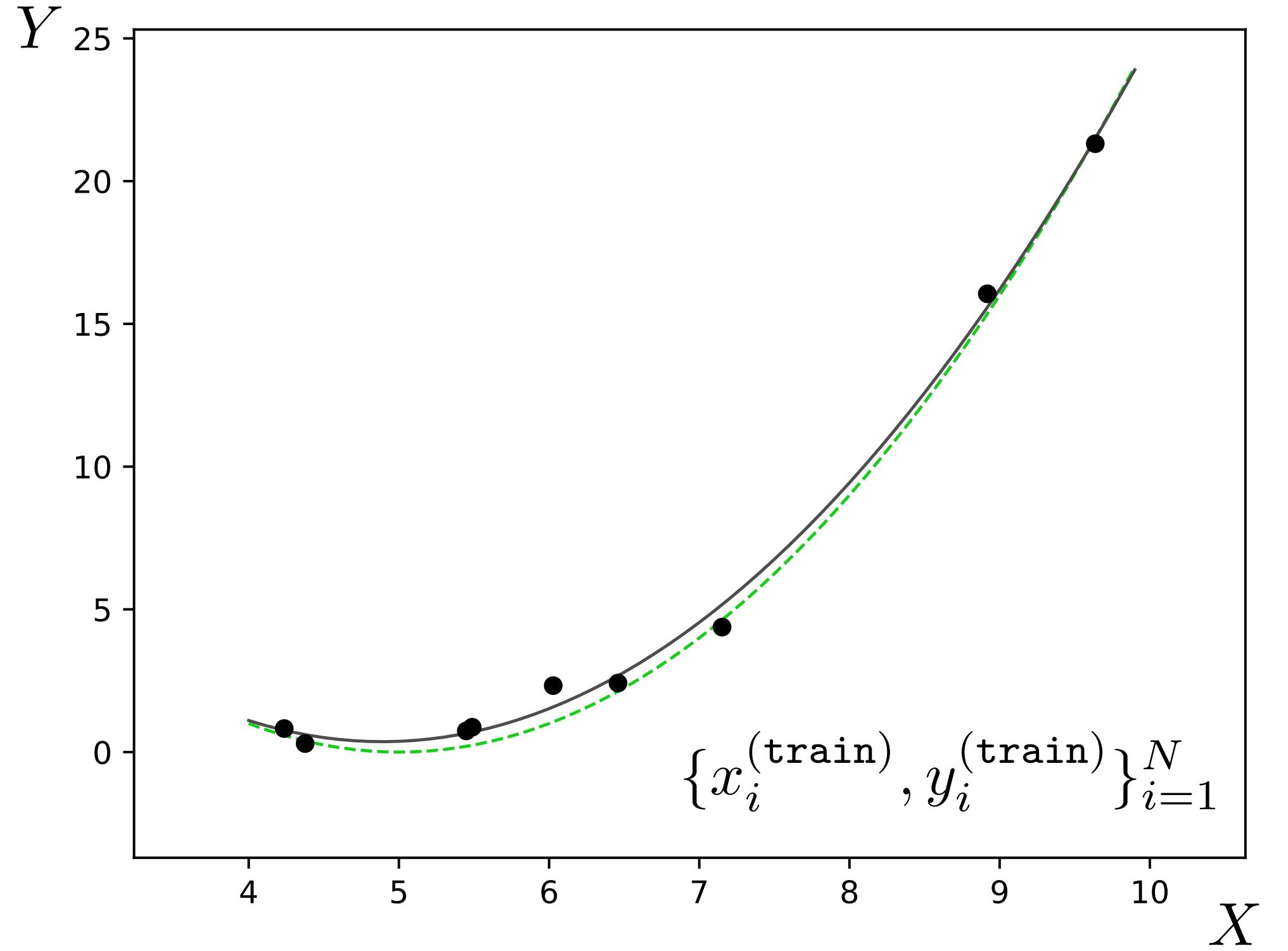
True data-generating process

$p_{\text{data}}$

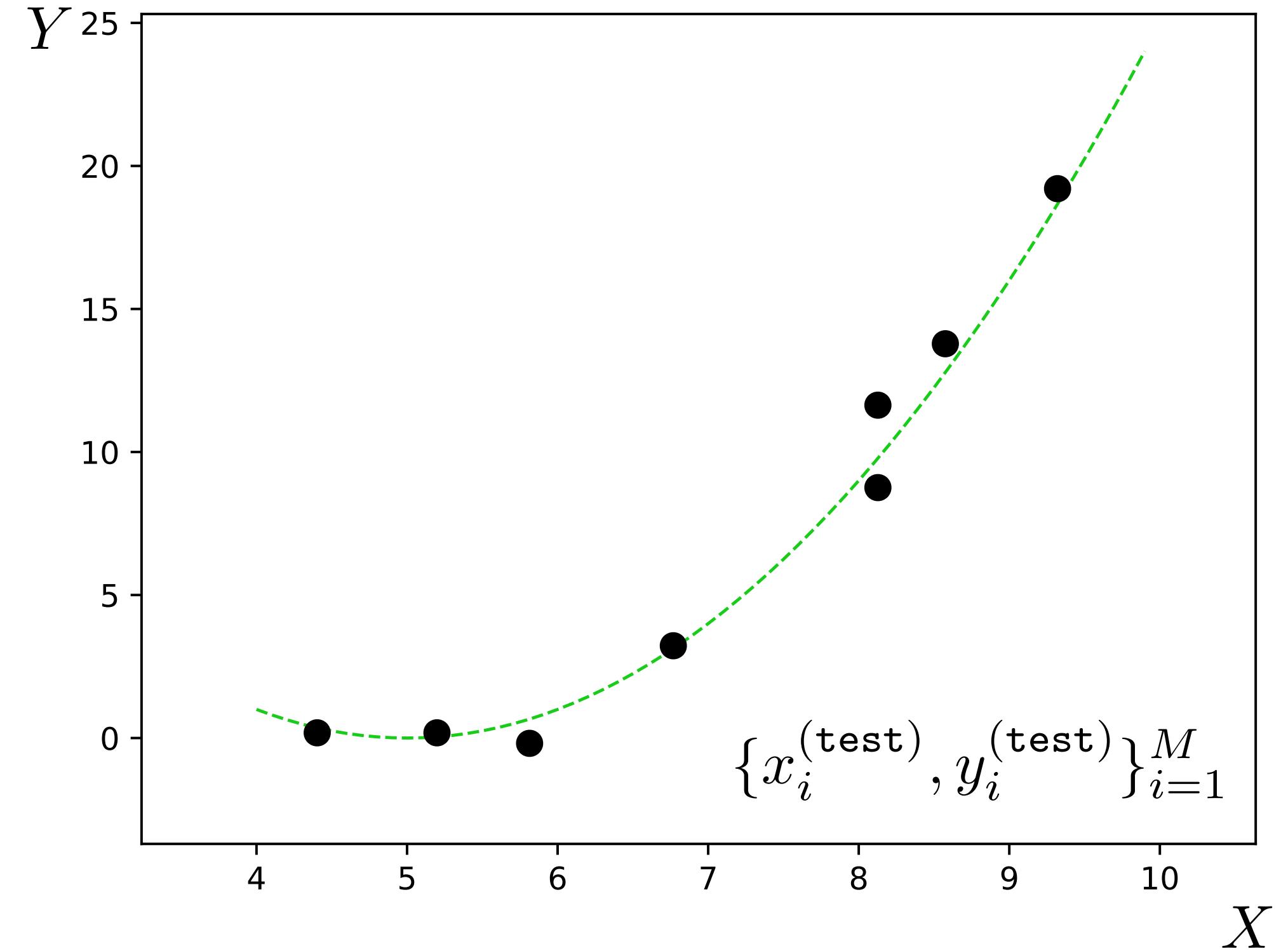
$$\{x_i^{(\text{train})}, y_i^{(\text{train})}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

$$\{x_i^{(\text{test})}, y_i^{(\text{test})}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

## Training data



## Test data



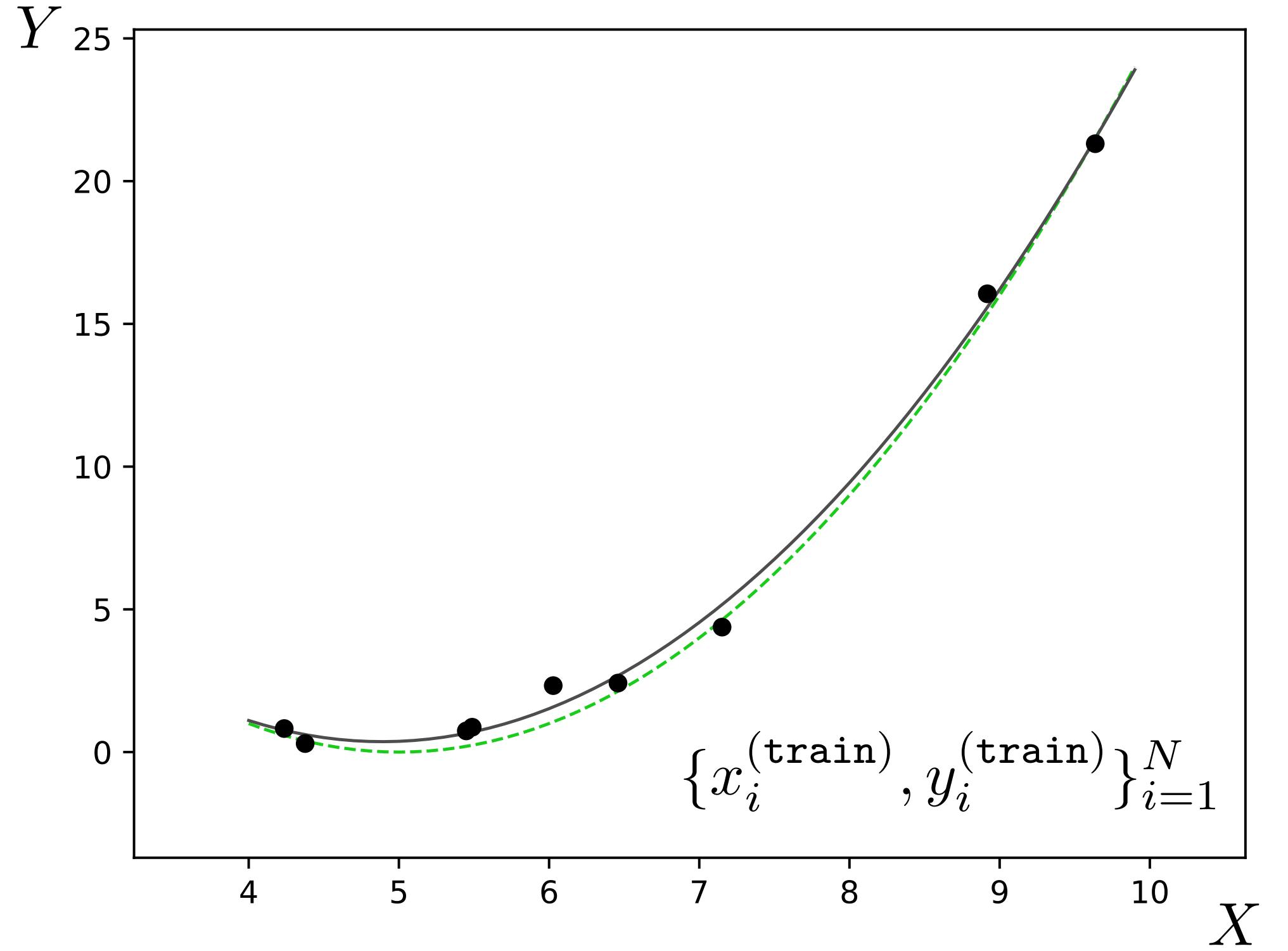
Training objective:

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

Test time evaluation:

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

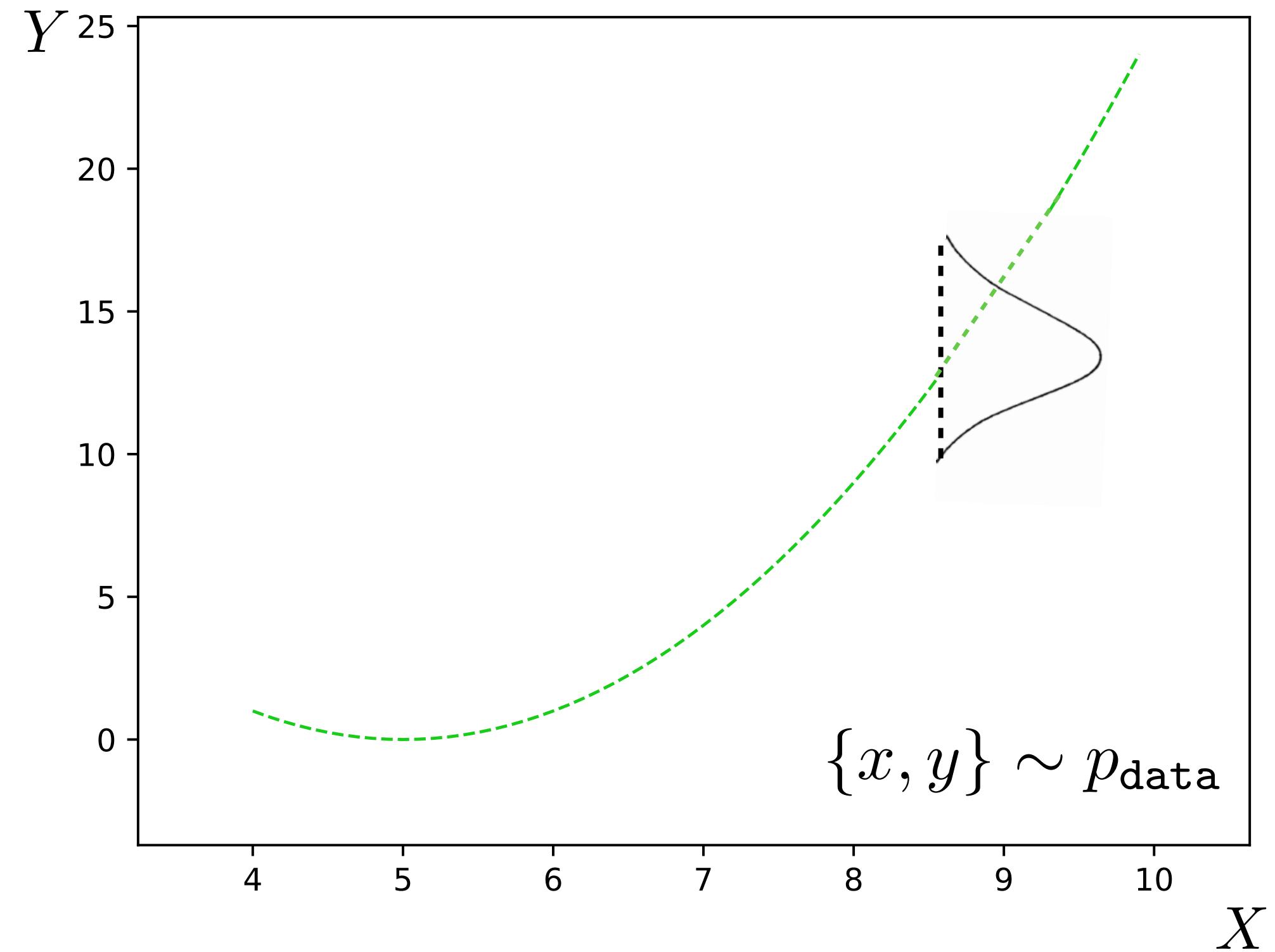
## Training data



Training objective:

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

## Test data



True objective:

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

# Generalization

“The central challenge in machine learning is that our algorithm must perform well on new, previously unseen inputs—not just those on which our model was trained. The ability to perform well on previously unobserved inputs is called **generalization**.

... [this is what] separates machine learning from optimization.”

— Deep Learning textbook (Goodfellow et al.)

# What does $\star$ do?

$$2 \star 3 = 36$$

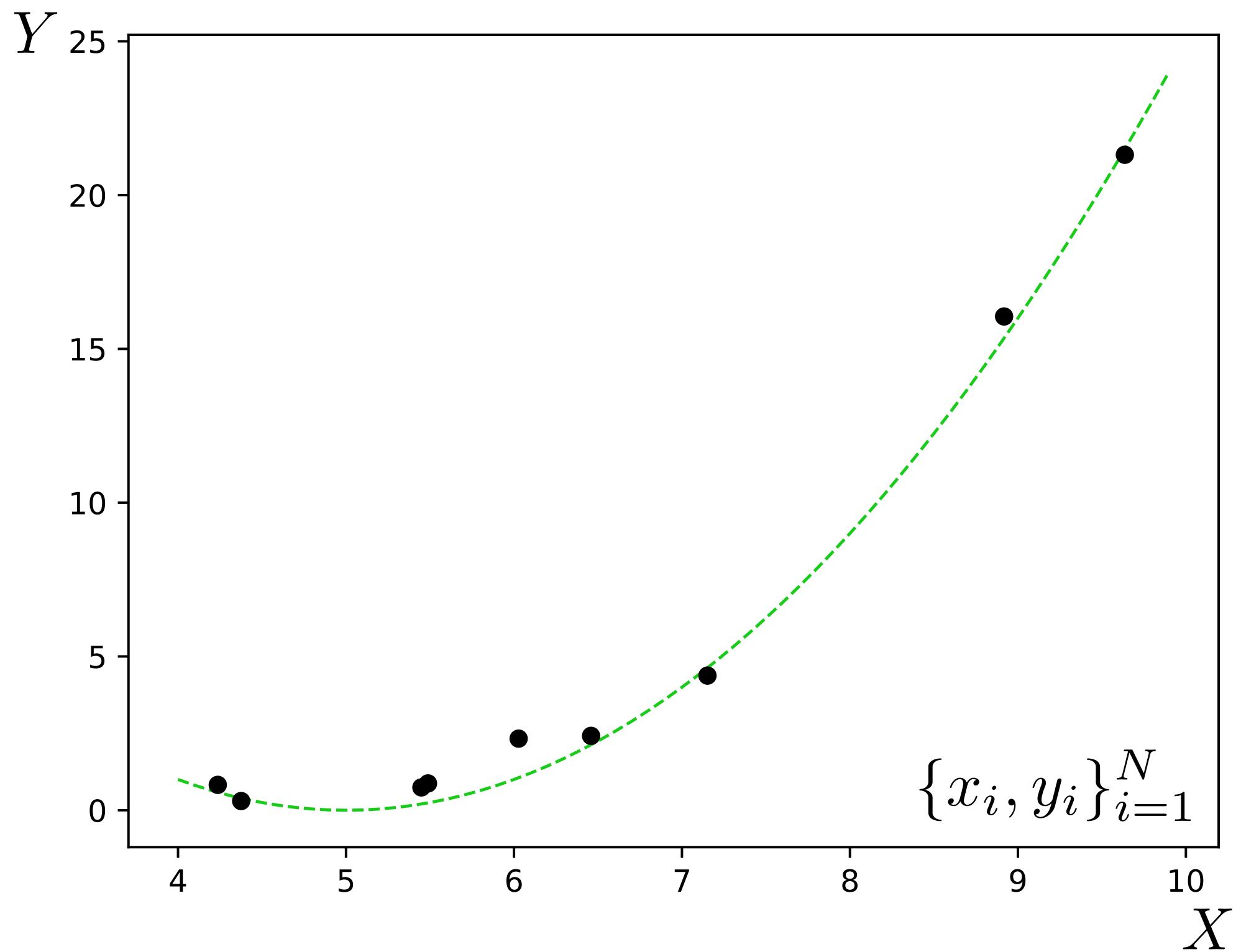
$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

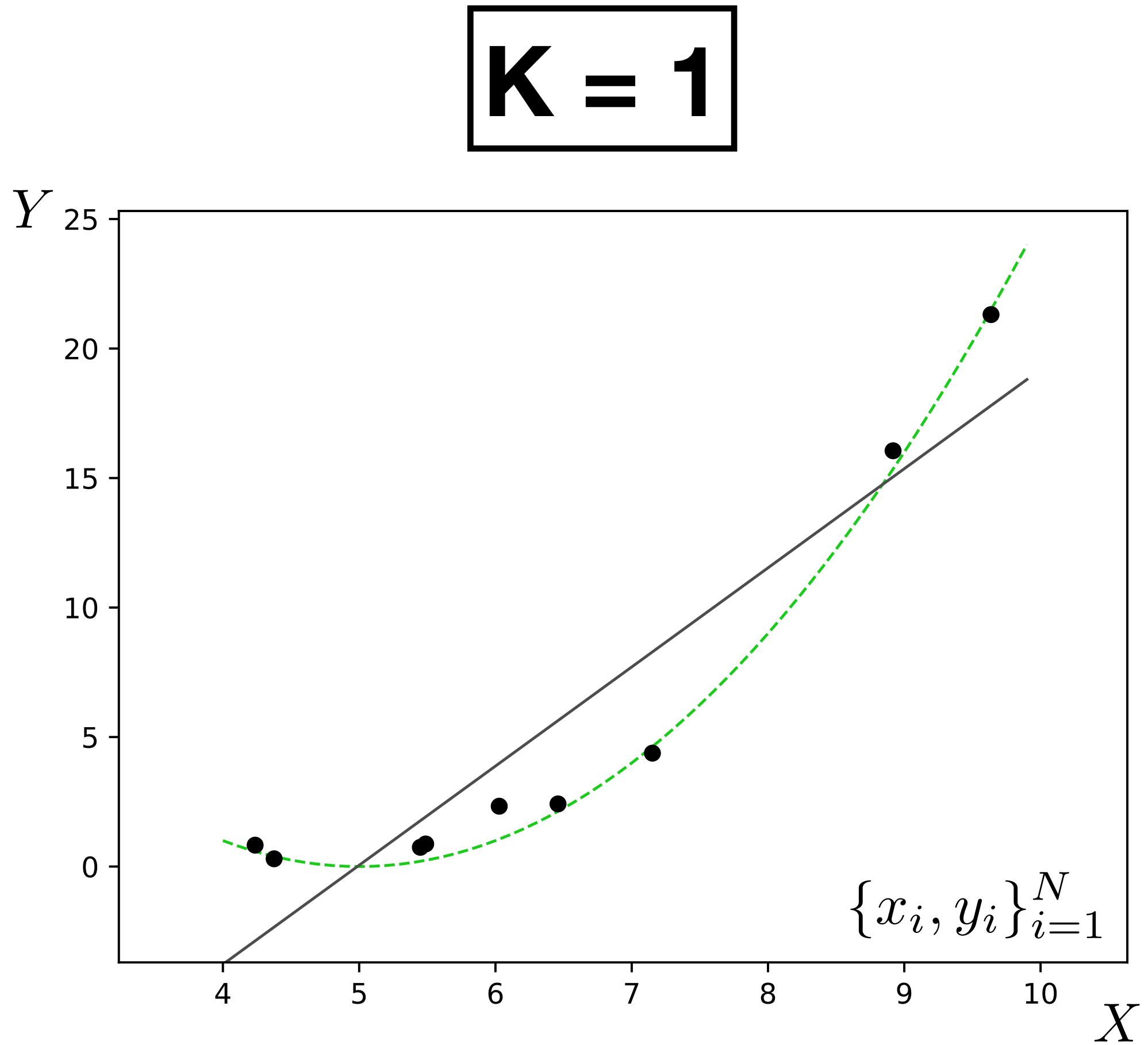
# What happens as we add more basis functions?

Training data



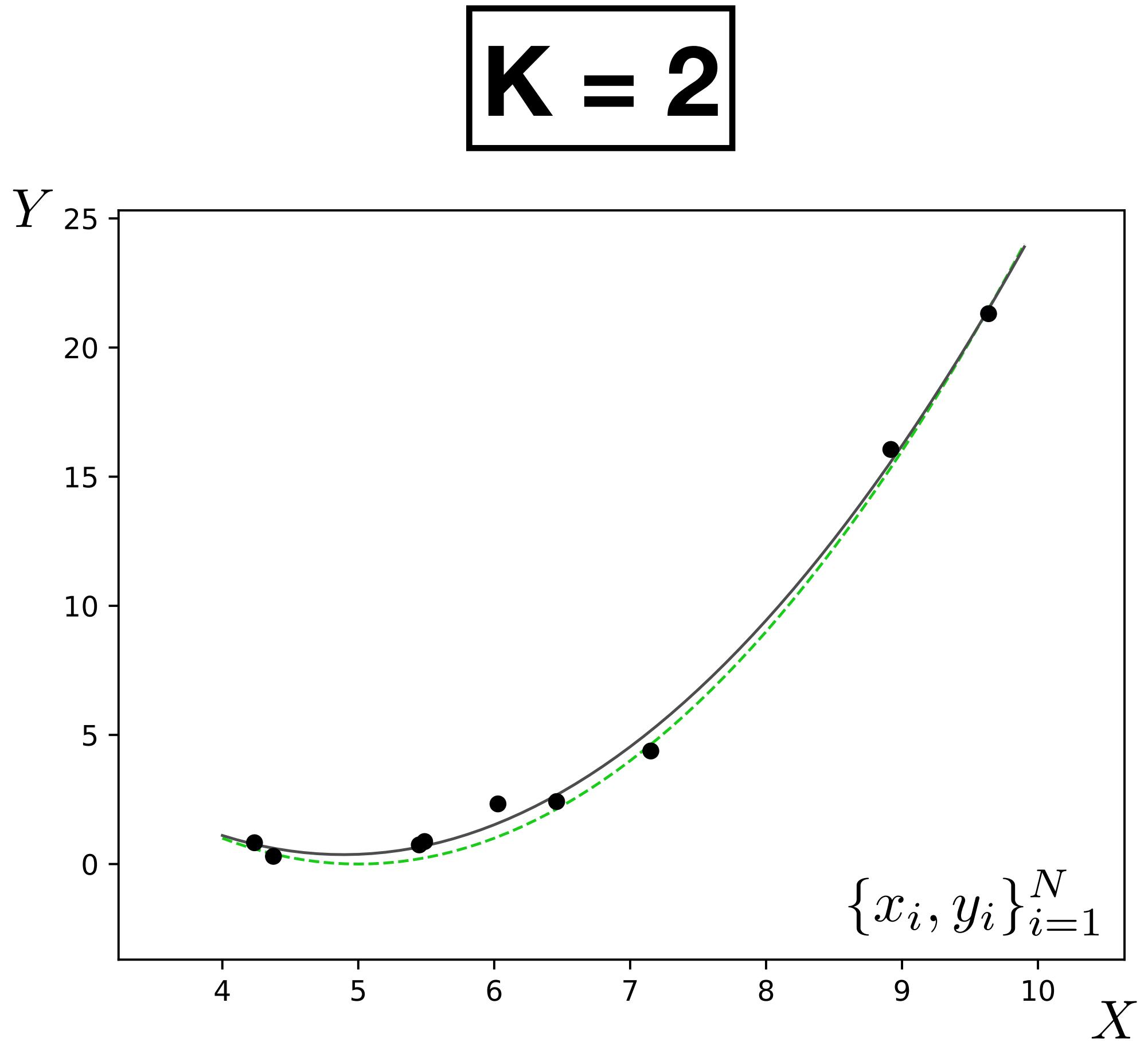
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



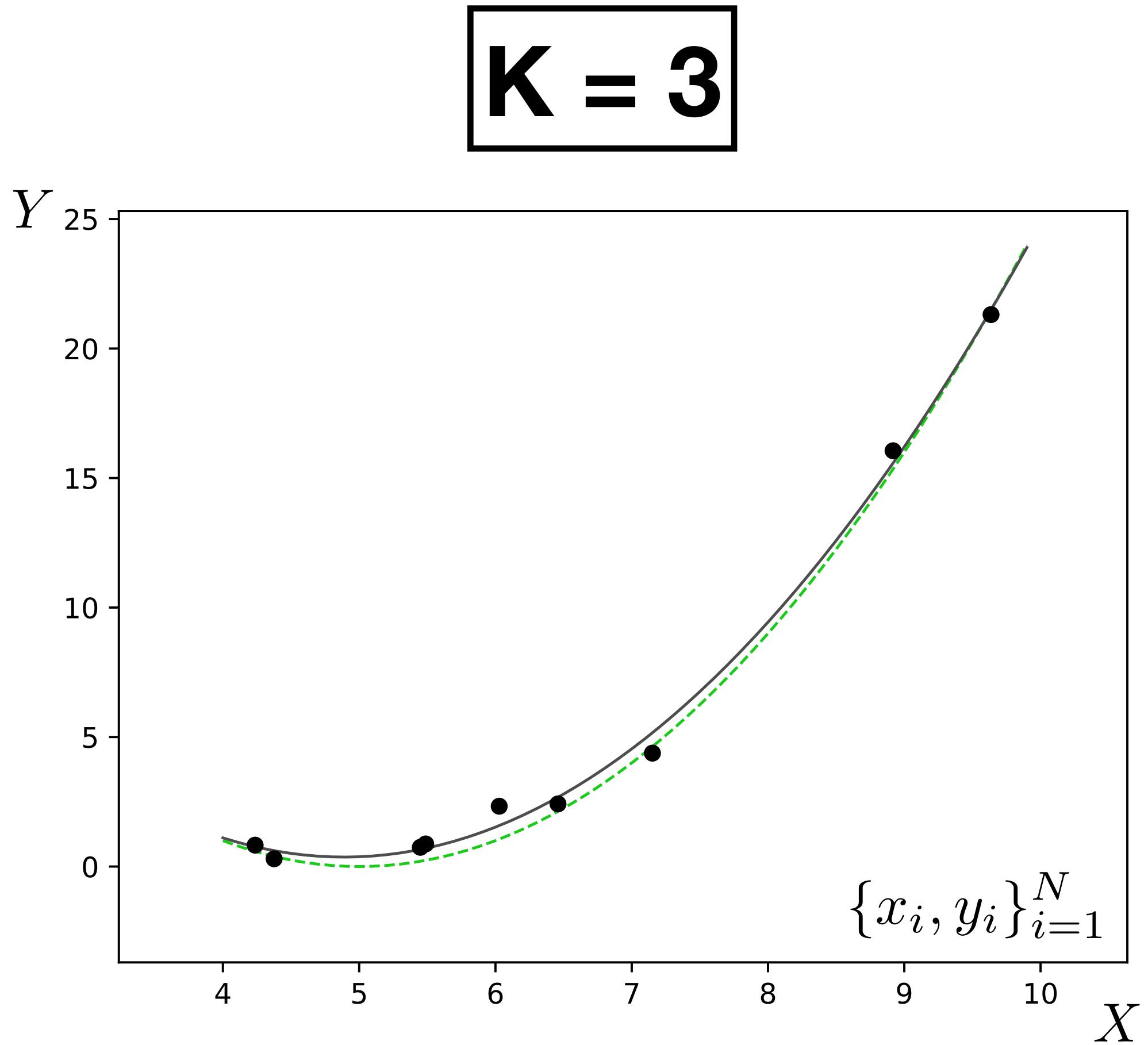
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



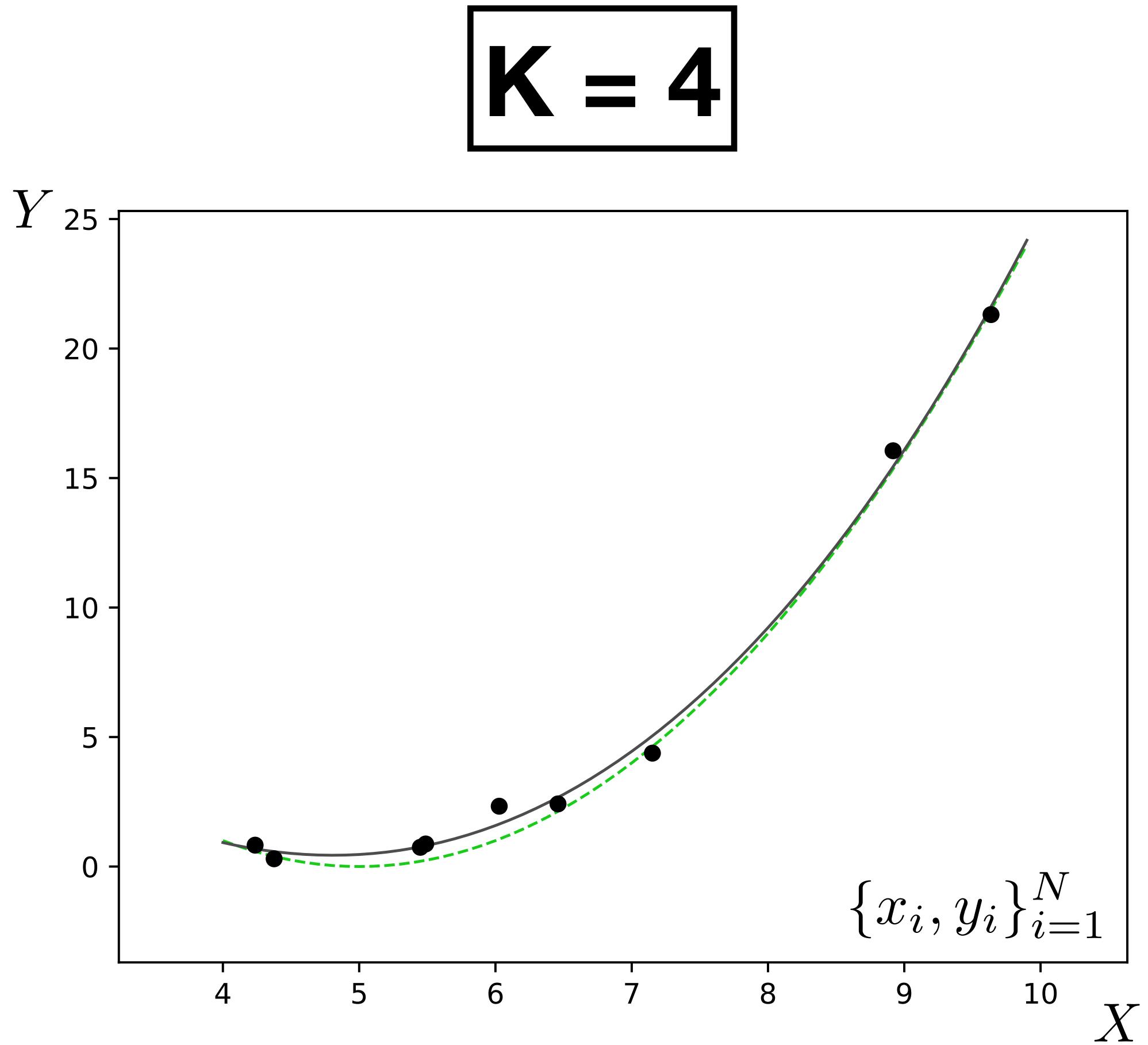
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



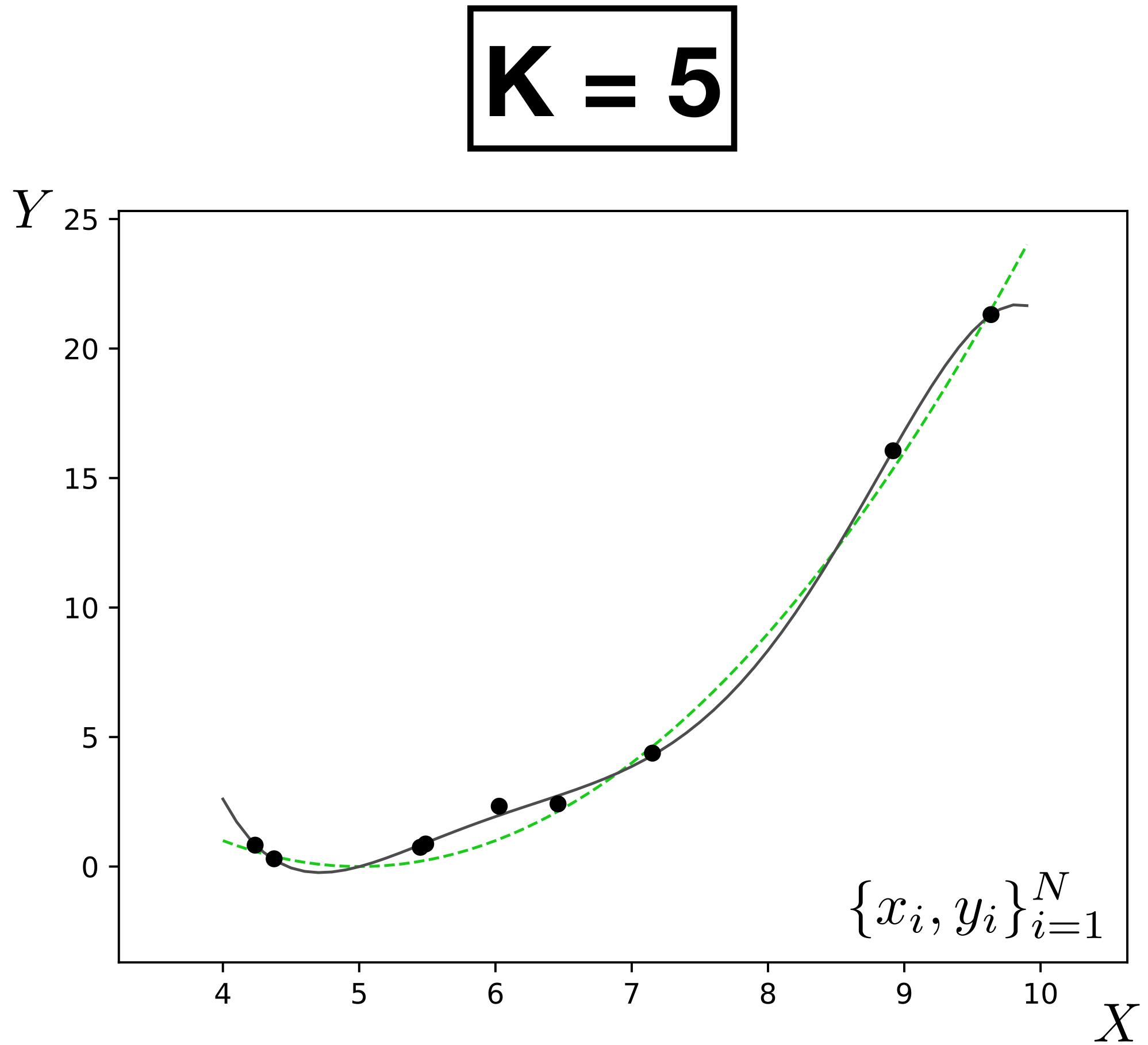
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



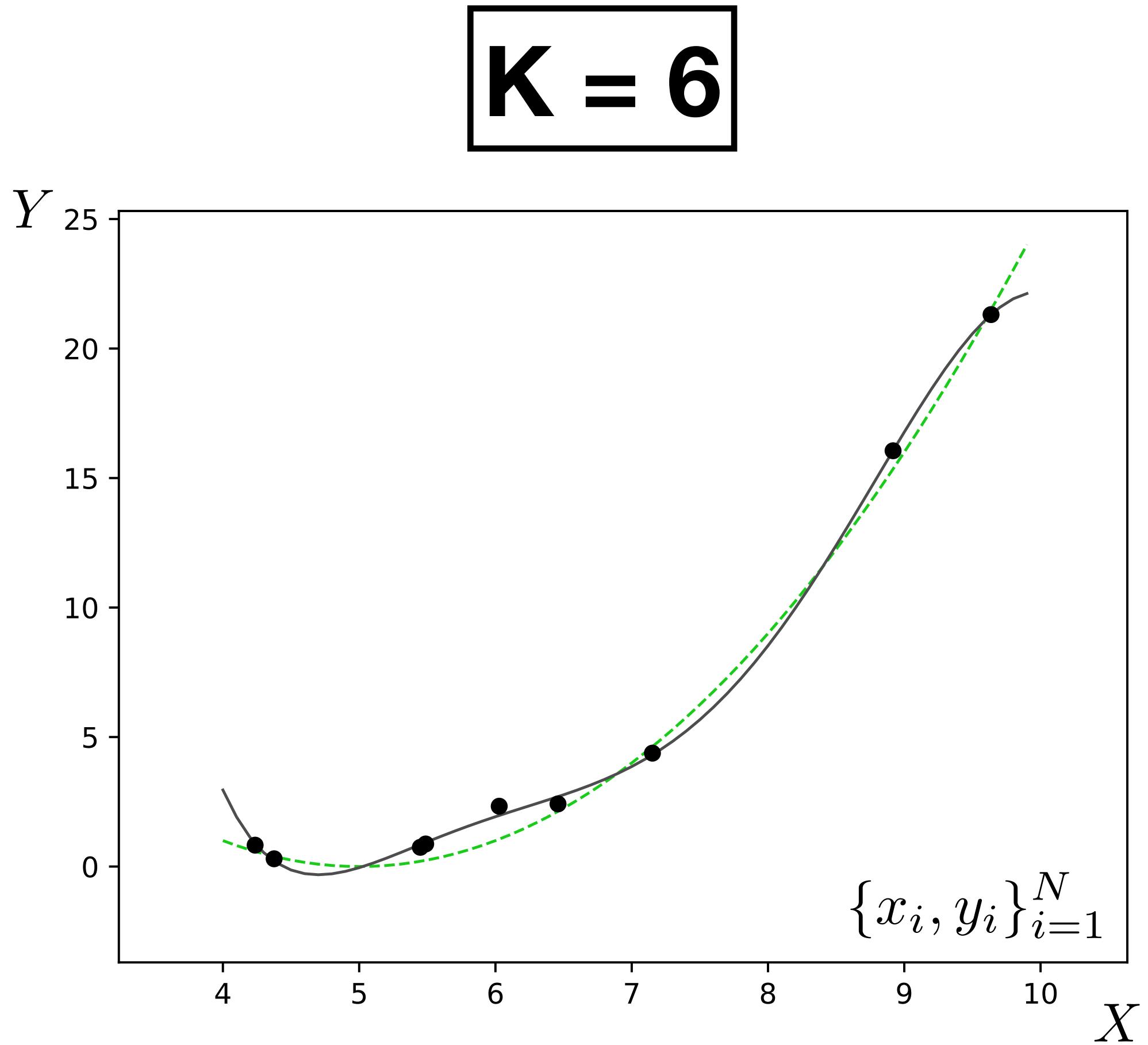
$$f_\theta(x) = \sum_{k=0}^K \theta_k x^k$$

# What happens as we add more basis functions?



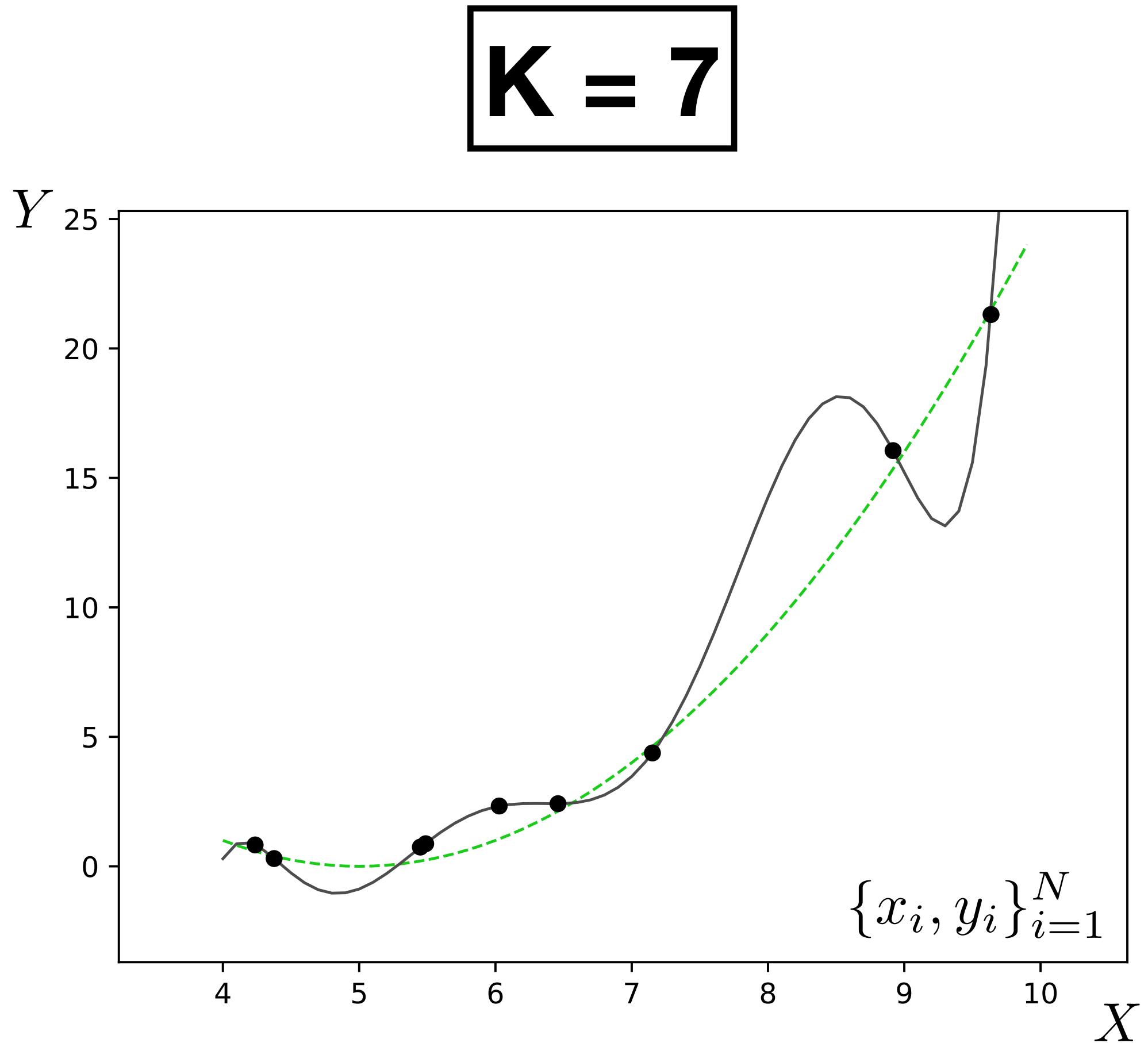
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



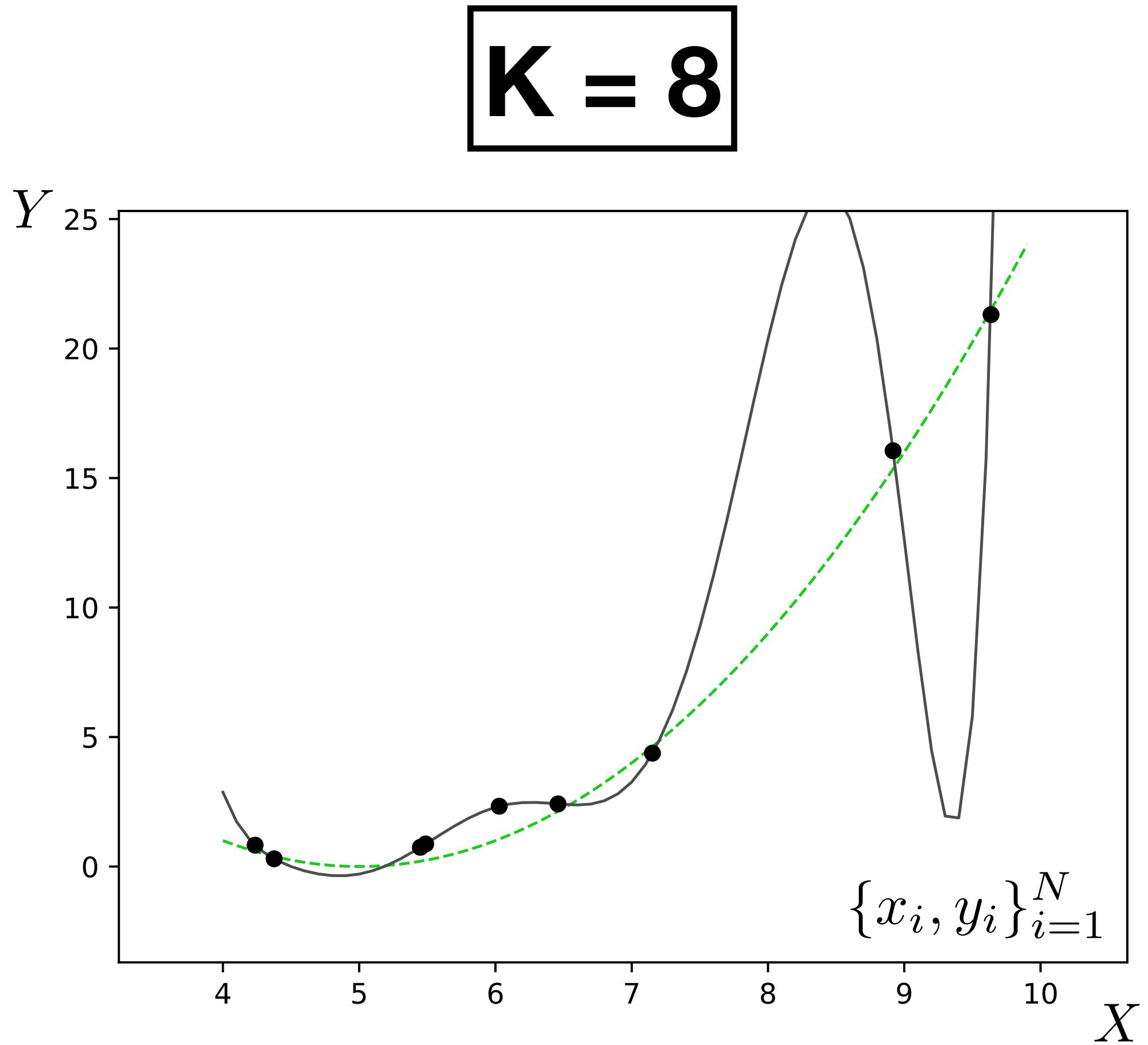
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



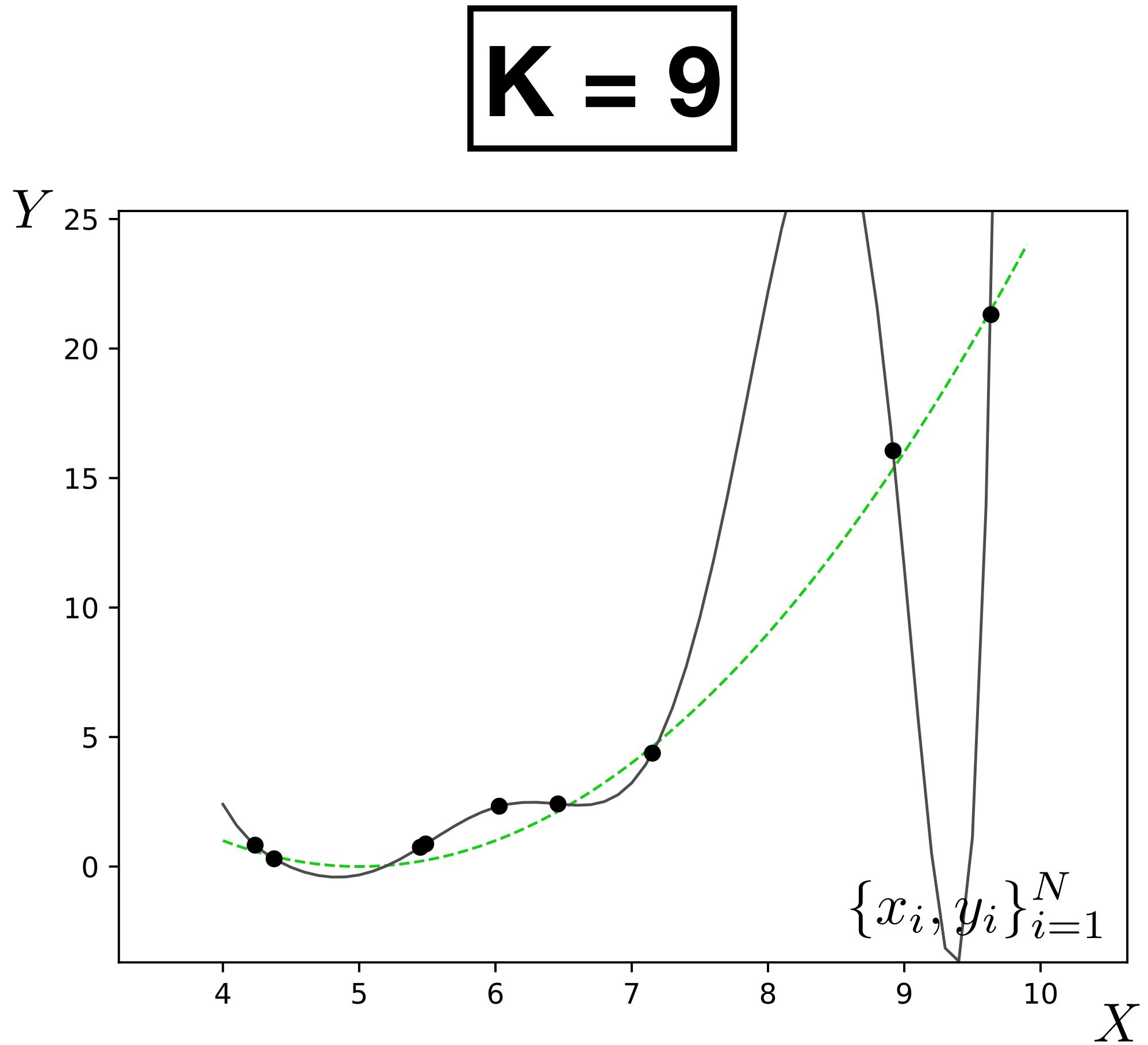
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



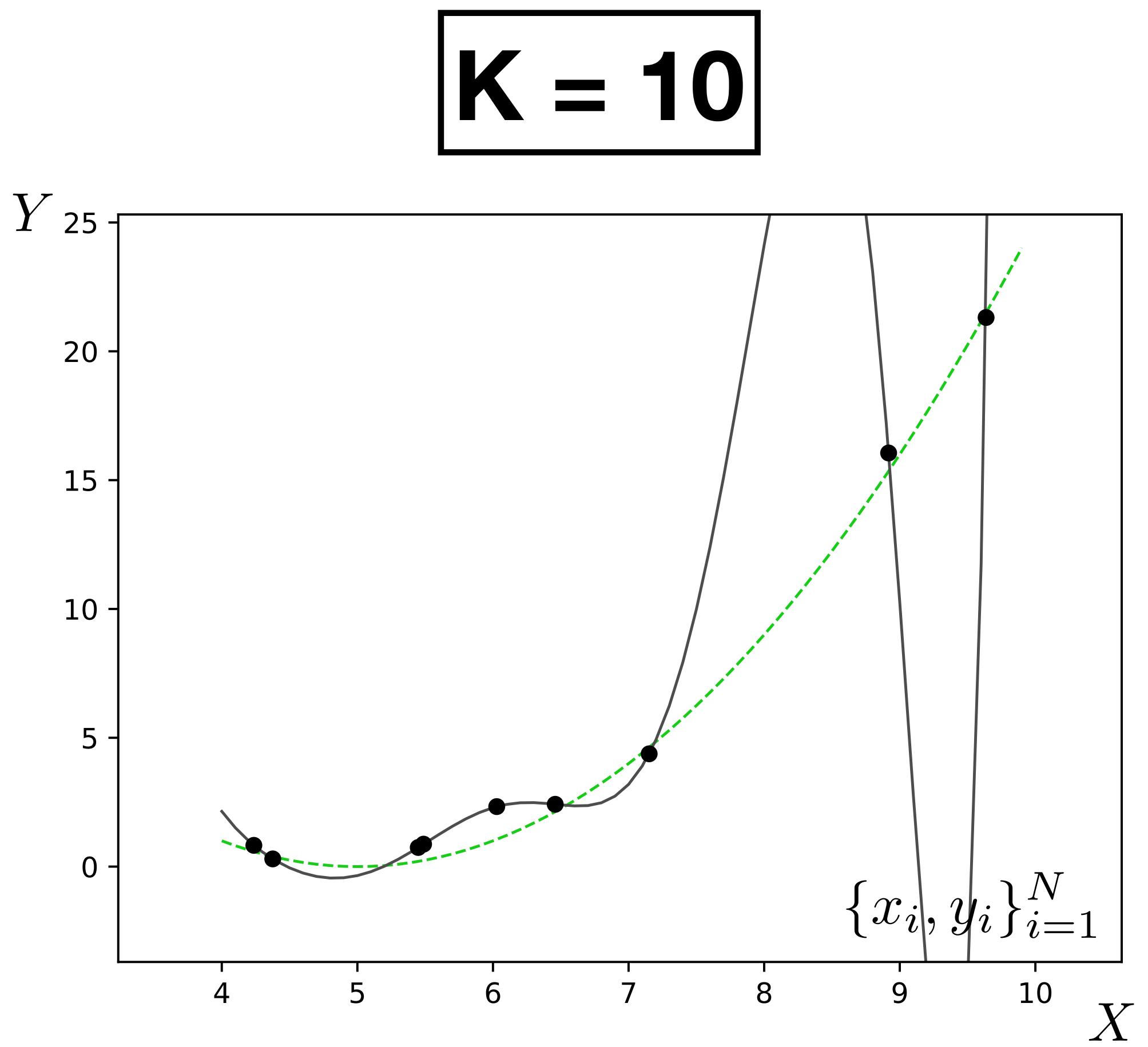
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

# What happens as we add more basis functions?



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

# What happens as we add more basis functions?

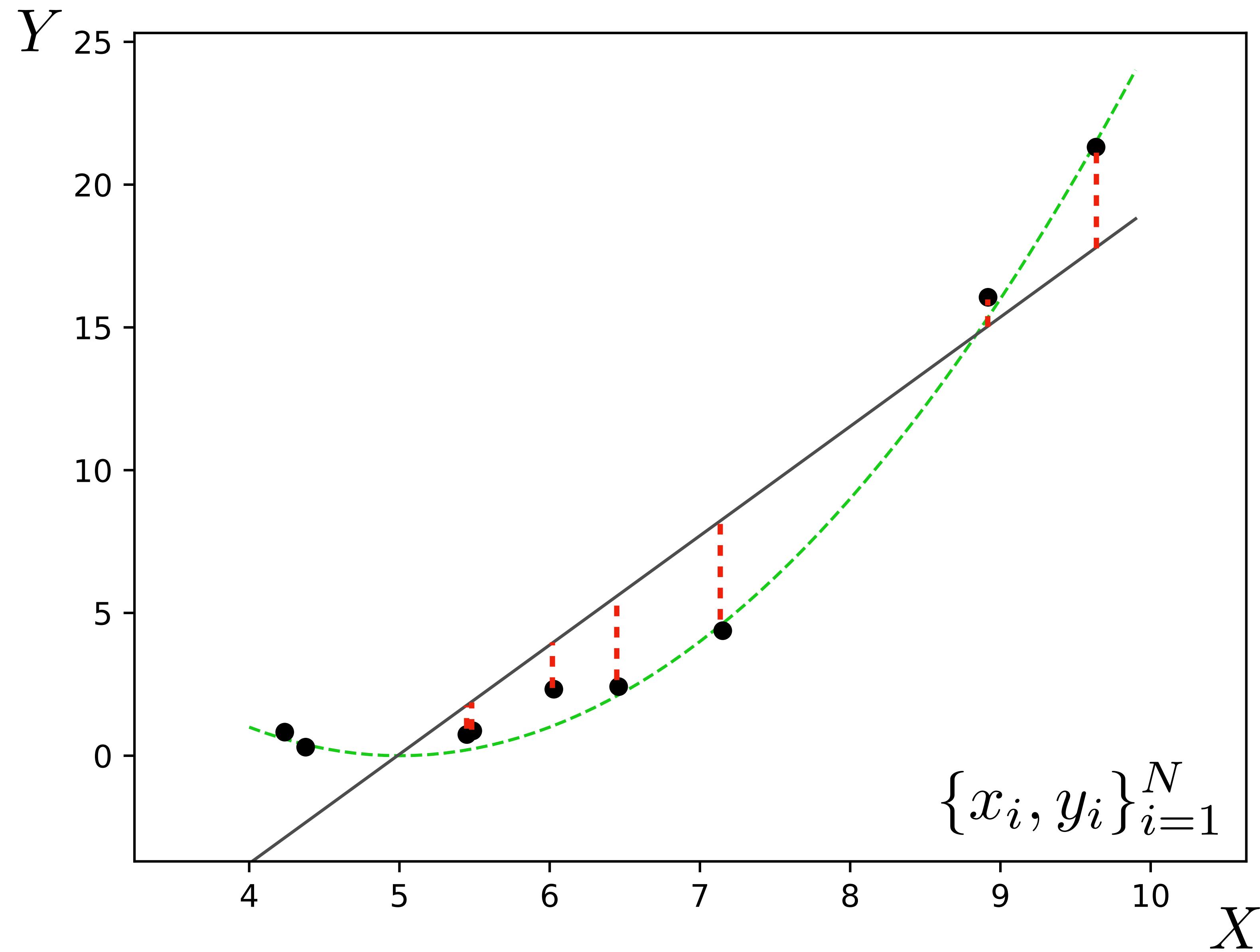


$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

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.

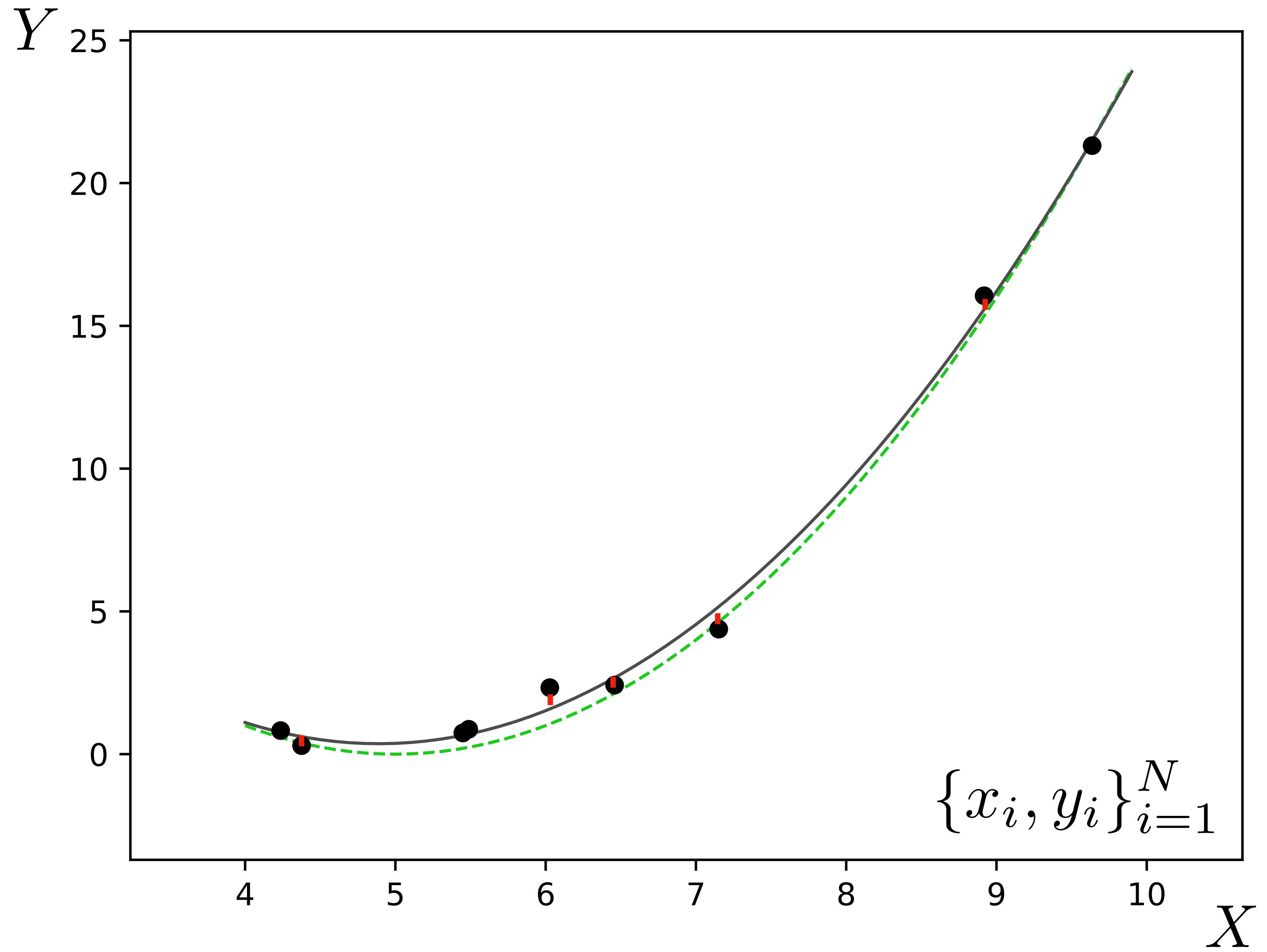
K = 1



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

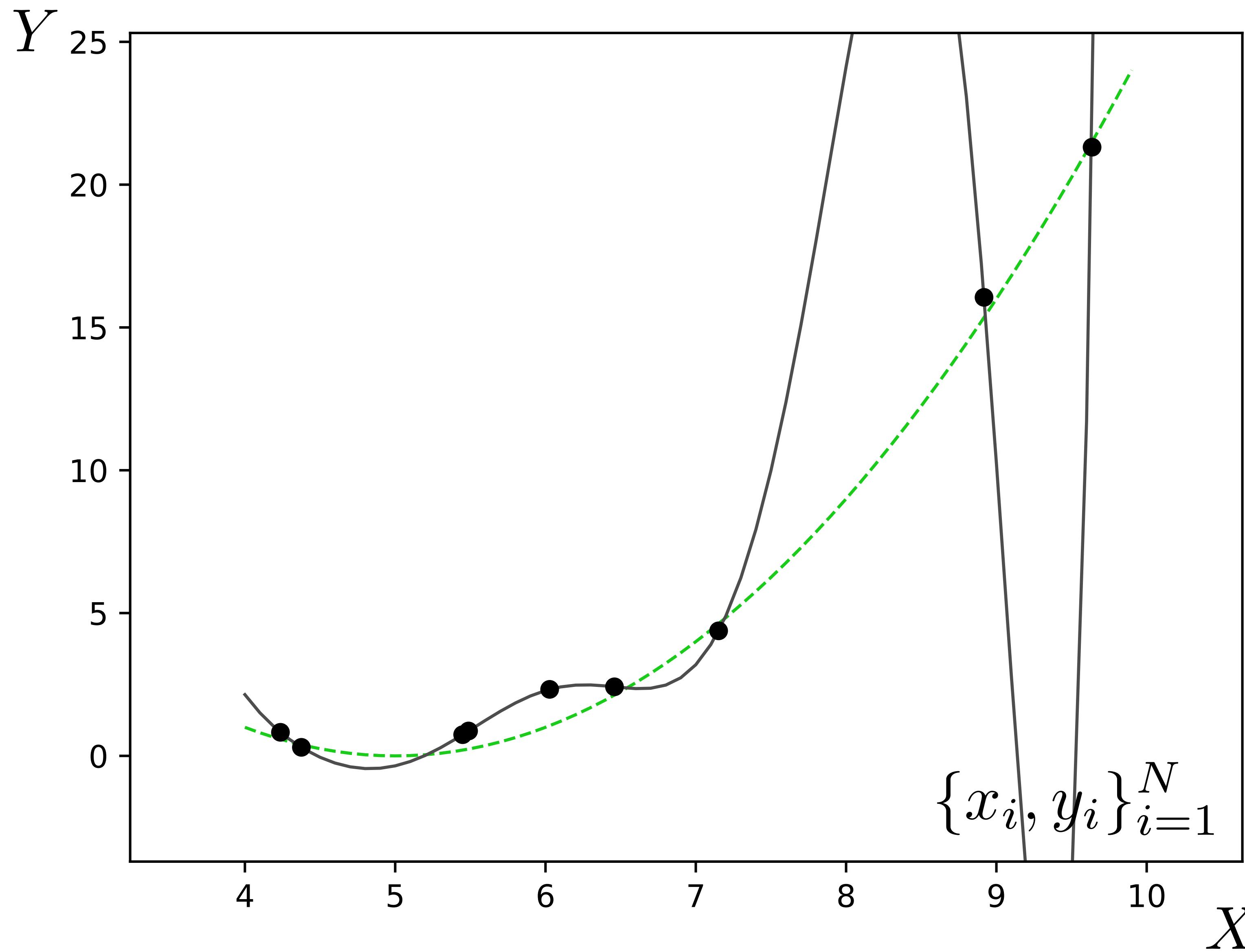
An underfit model will have high **error** on the training points. This error is known as **approximation error**.

K = 2



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

K = 10

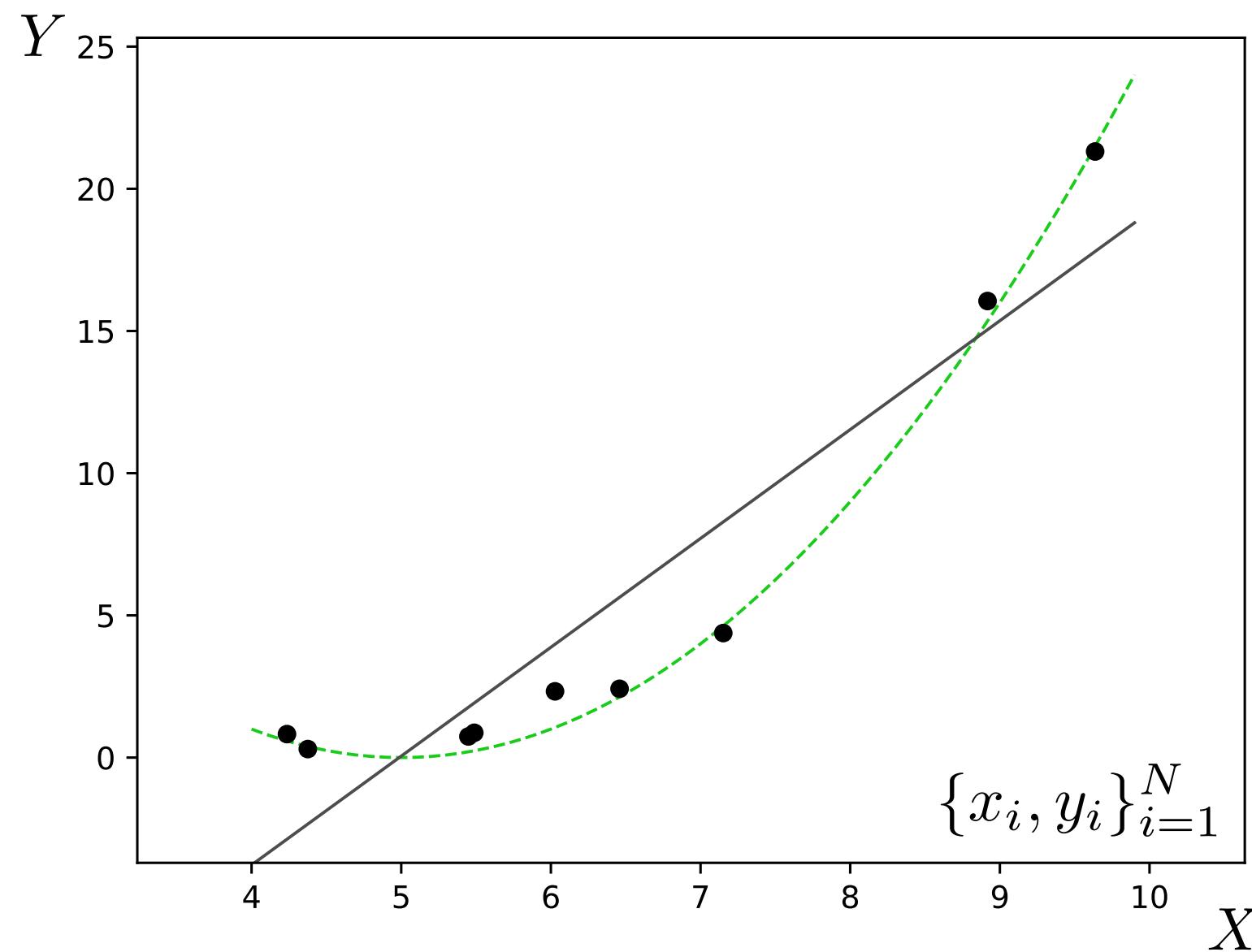


Now we have zero approximation error – the curve passes exactly through each training point.

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).

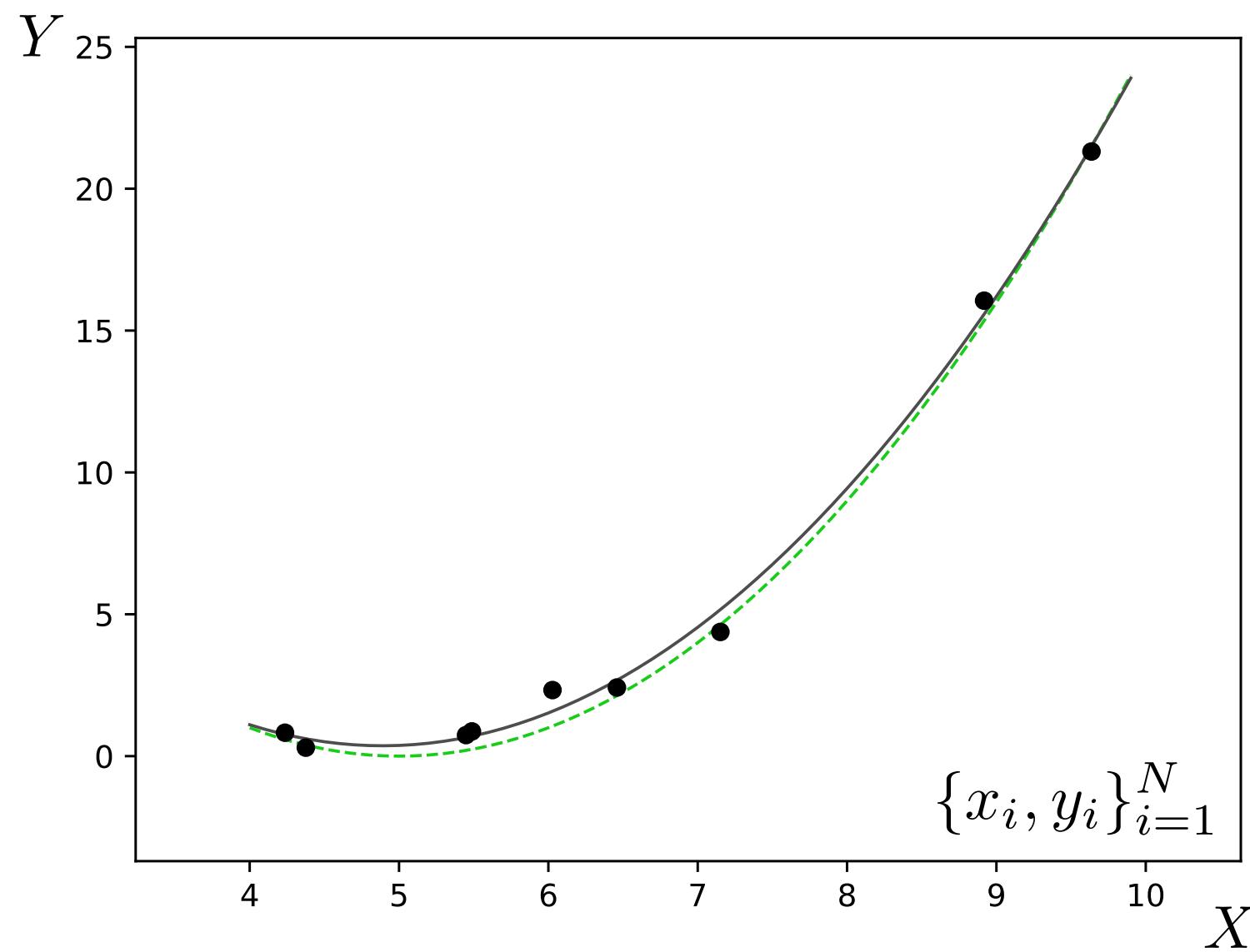
Underfitting

$$K = 1$$



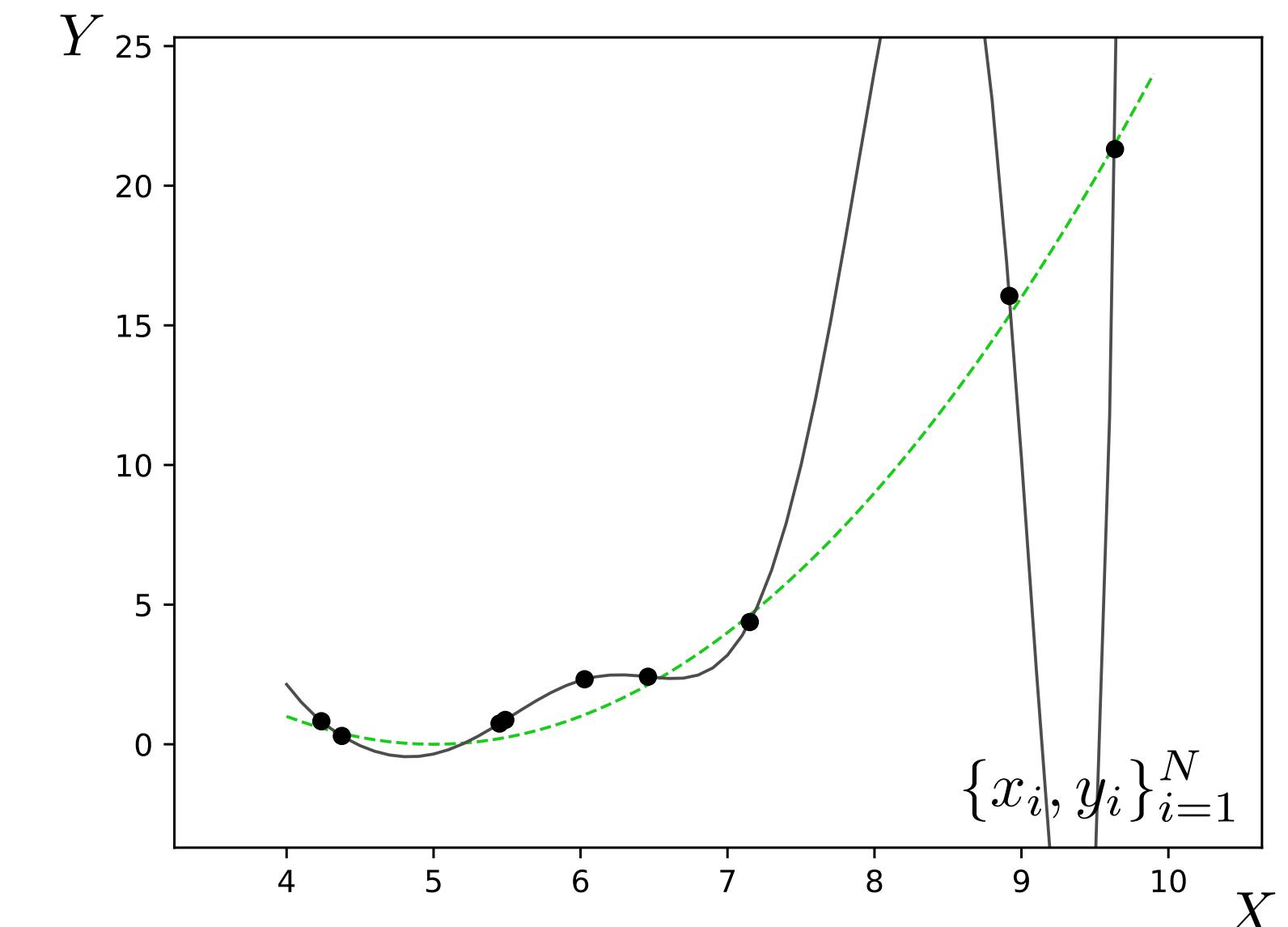
Appropriate model

$$K = 2$$



Overfitting

$$K = 10$$



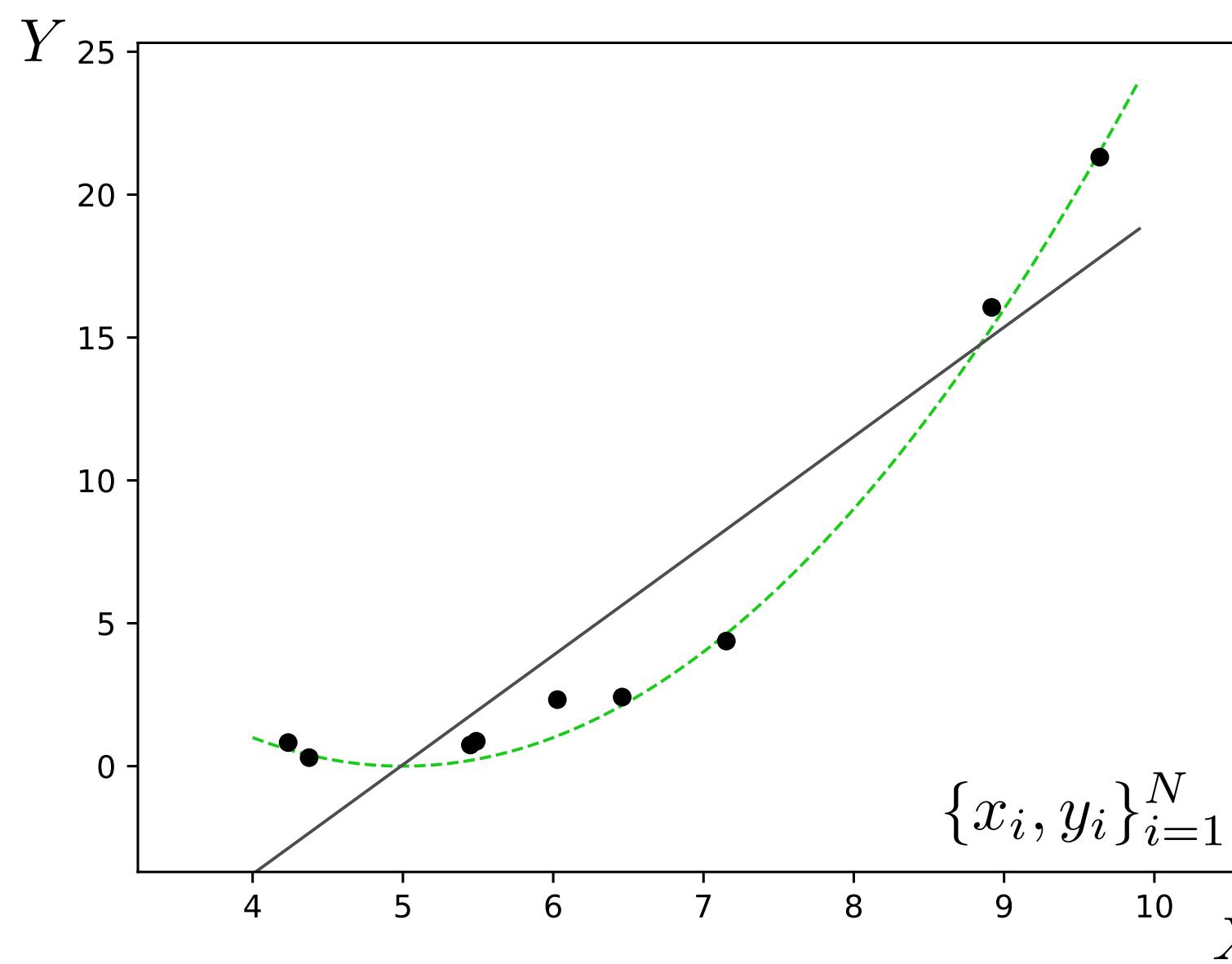
High error on train set  
High error on test set

Low error on train set  
Low error on test set

Lowest error on train set  
High error on test set

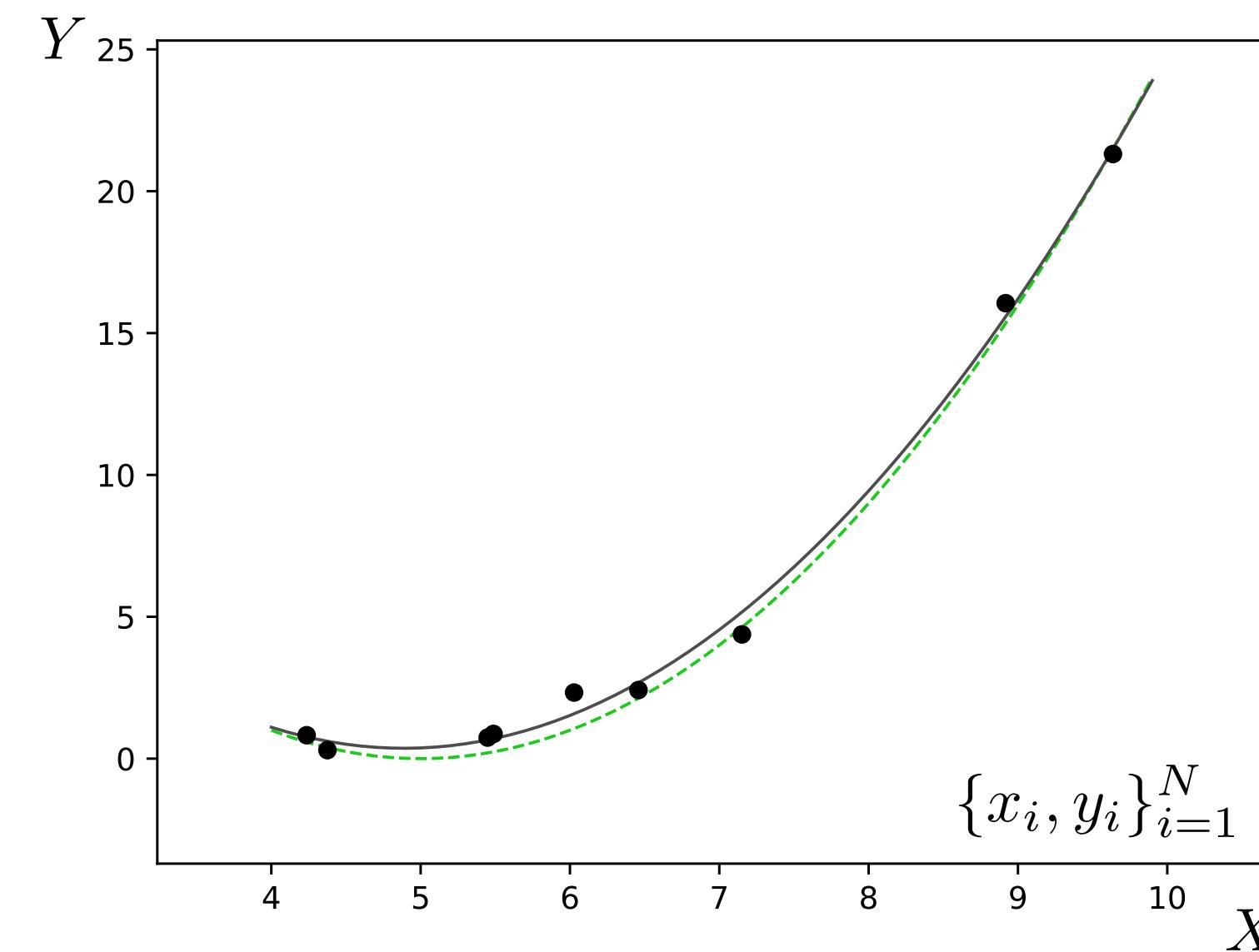
Underfitting

$K = 1$



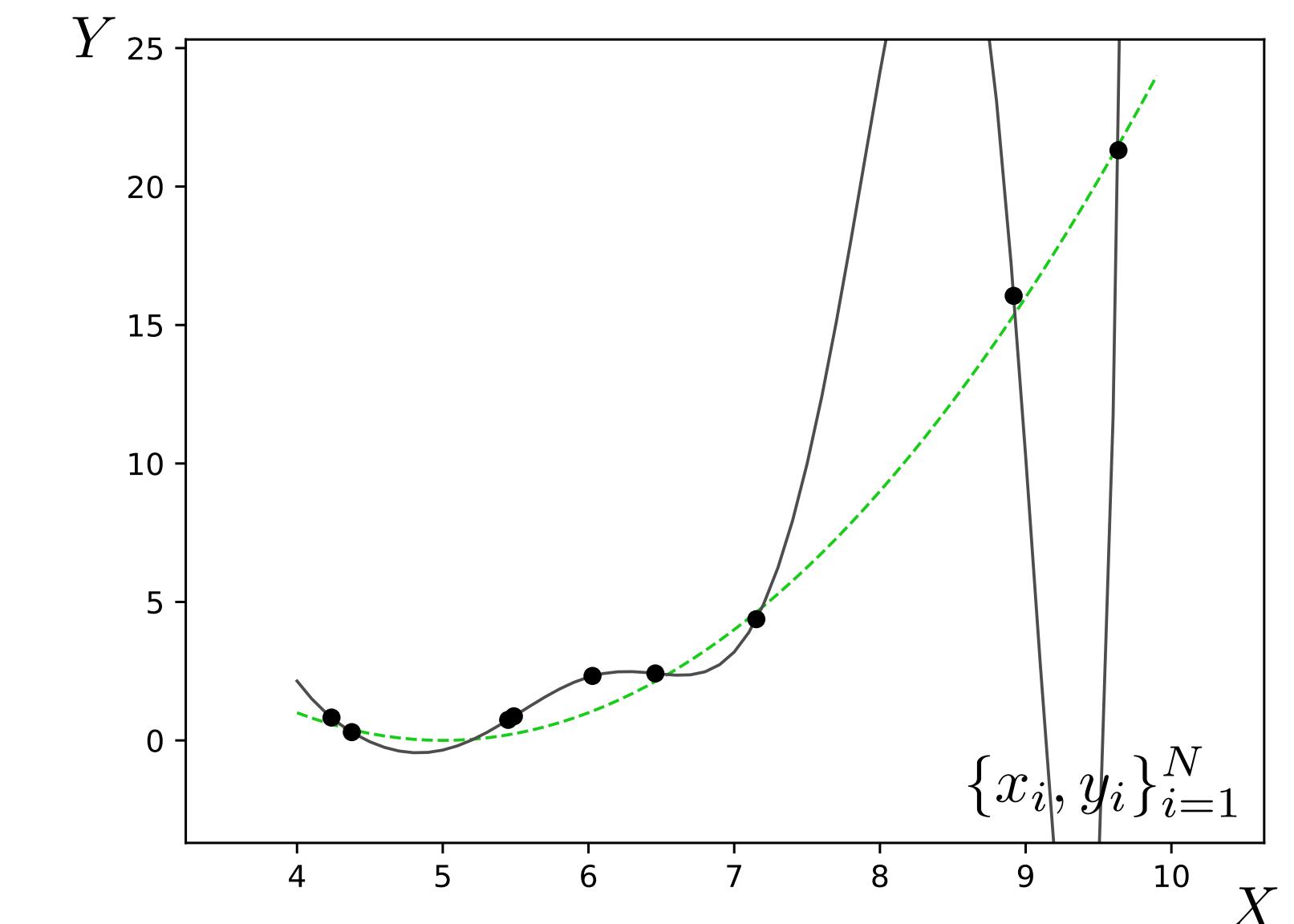
Appropriate model

$K = 2$



Overfitting

$K = 10$



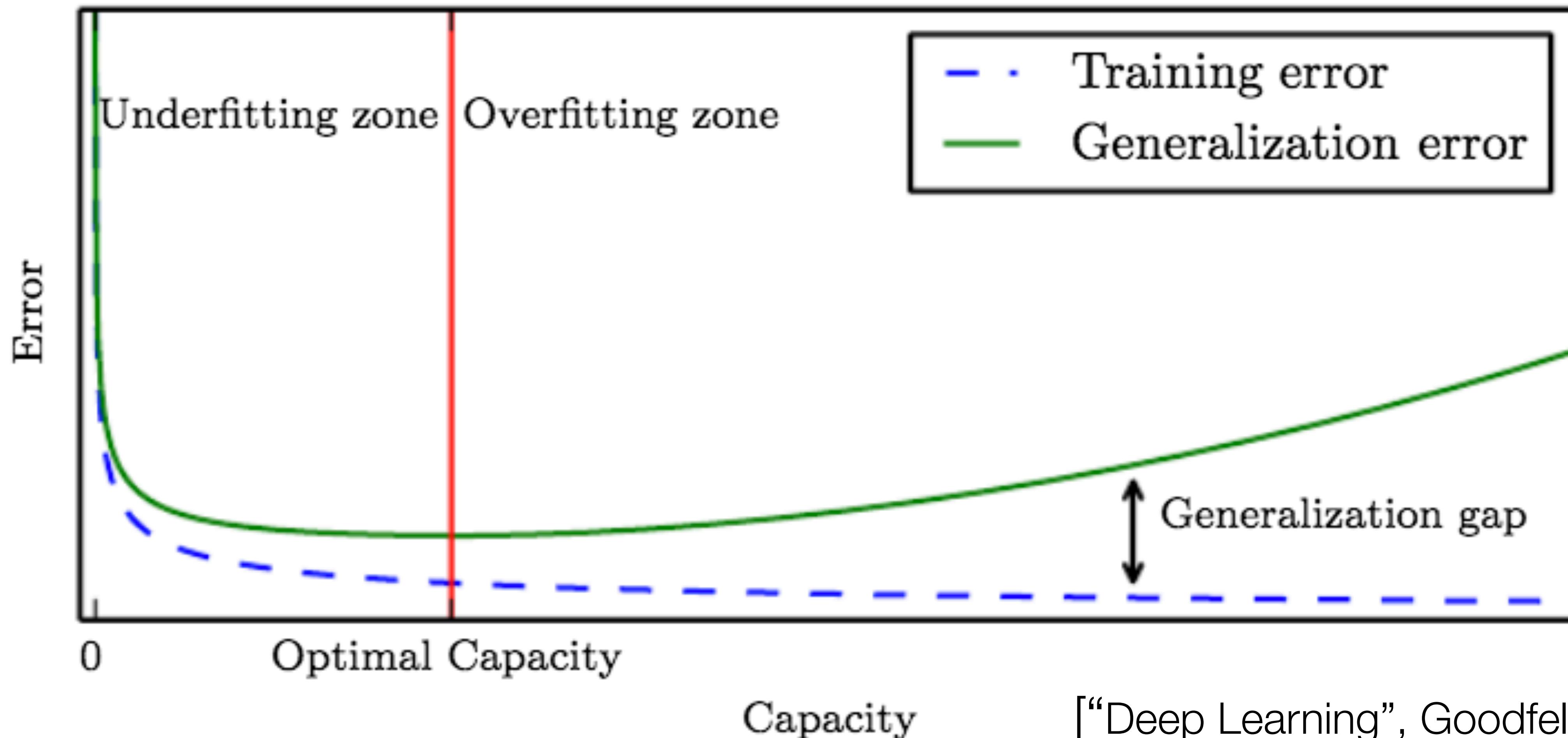
We need to control the **capacity** of the model (e.g., use the appropriate number of free parameters).

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.

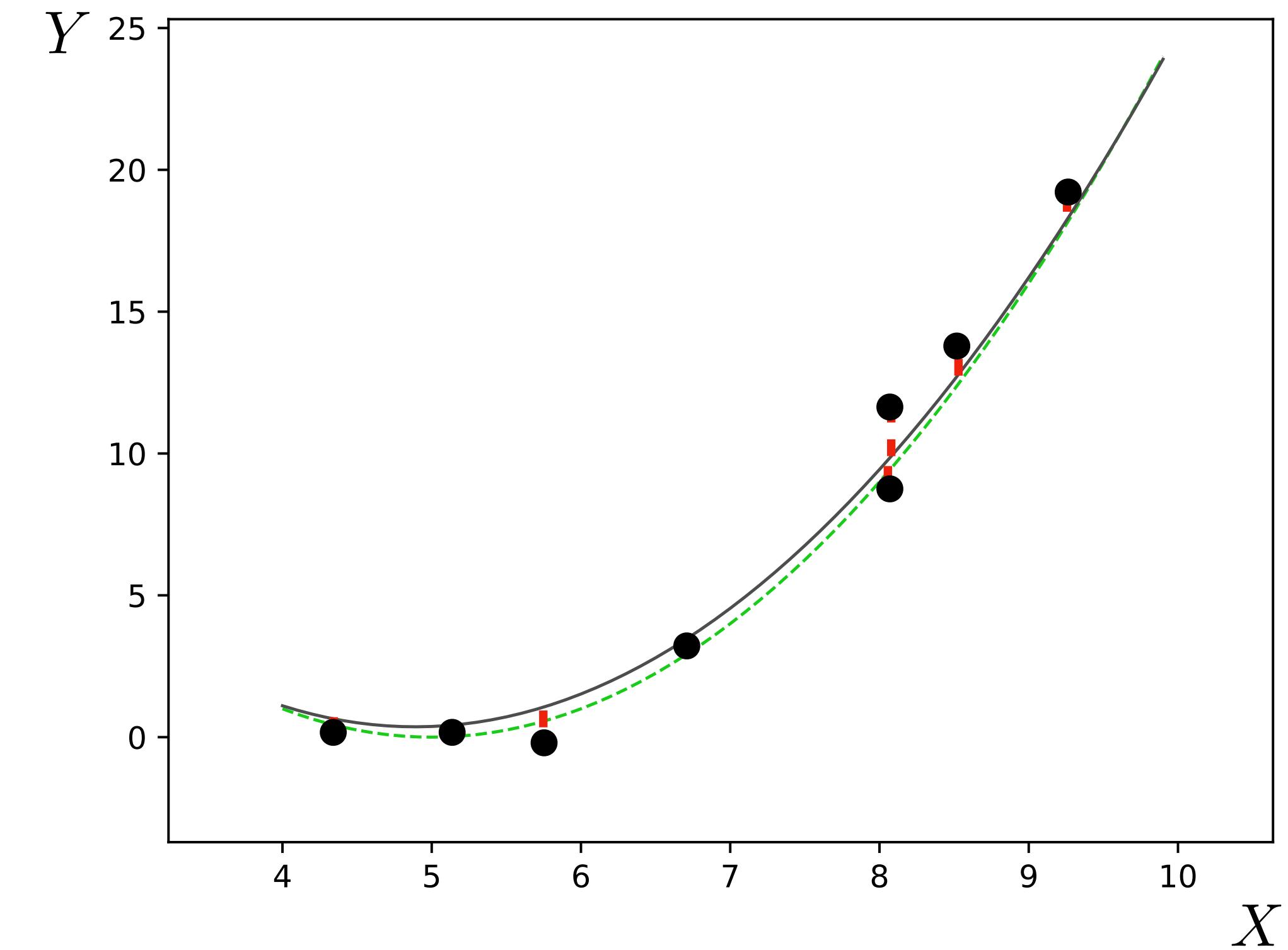
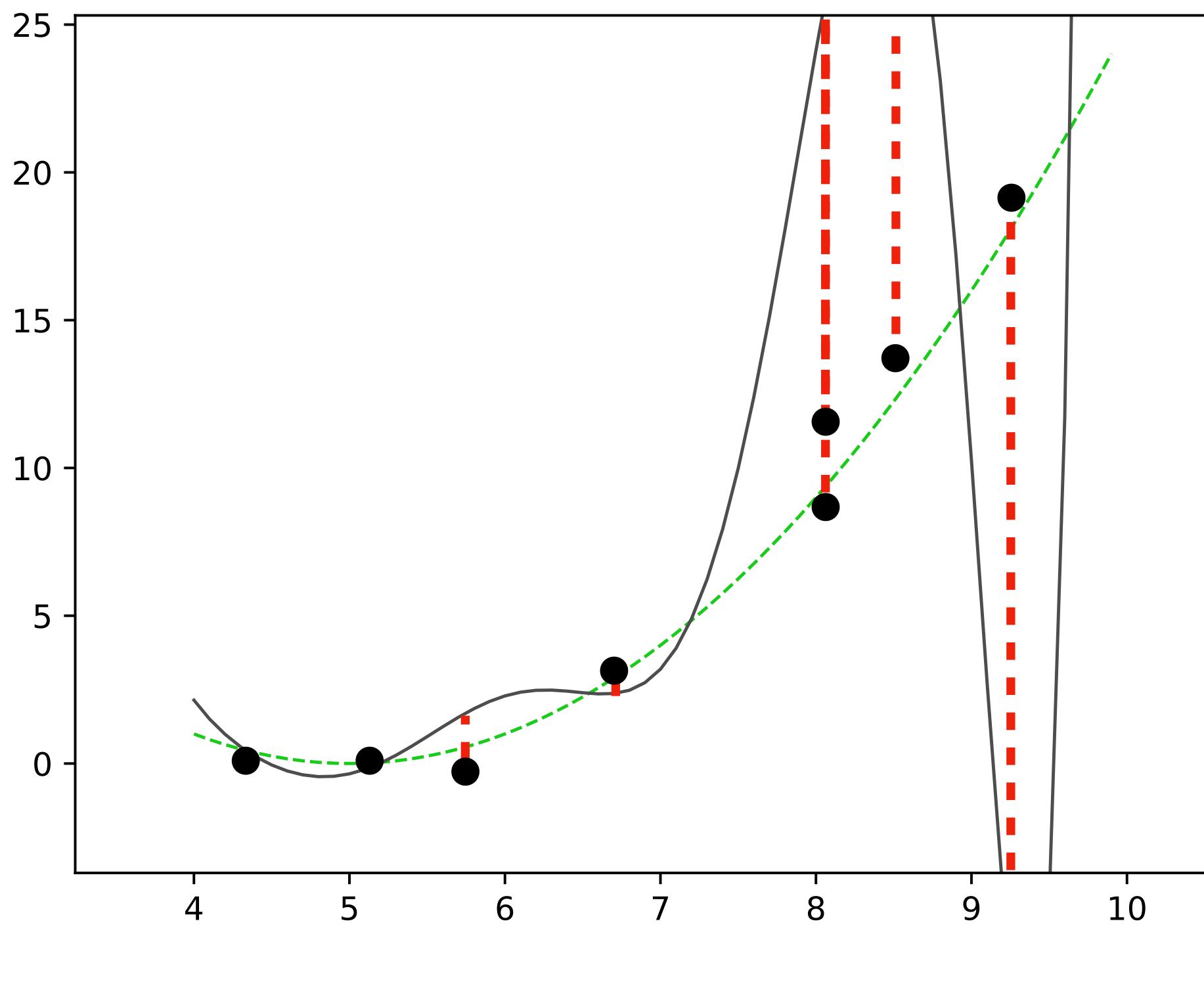
Simple models have low capacity.

# Training error versus generalization error



# 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 data

# Fitting just right



Underfitting?

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

Overfitting?

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

Selecting a *hypothesis* space of functions with just the right capacity is known as **model selection**

# 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

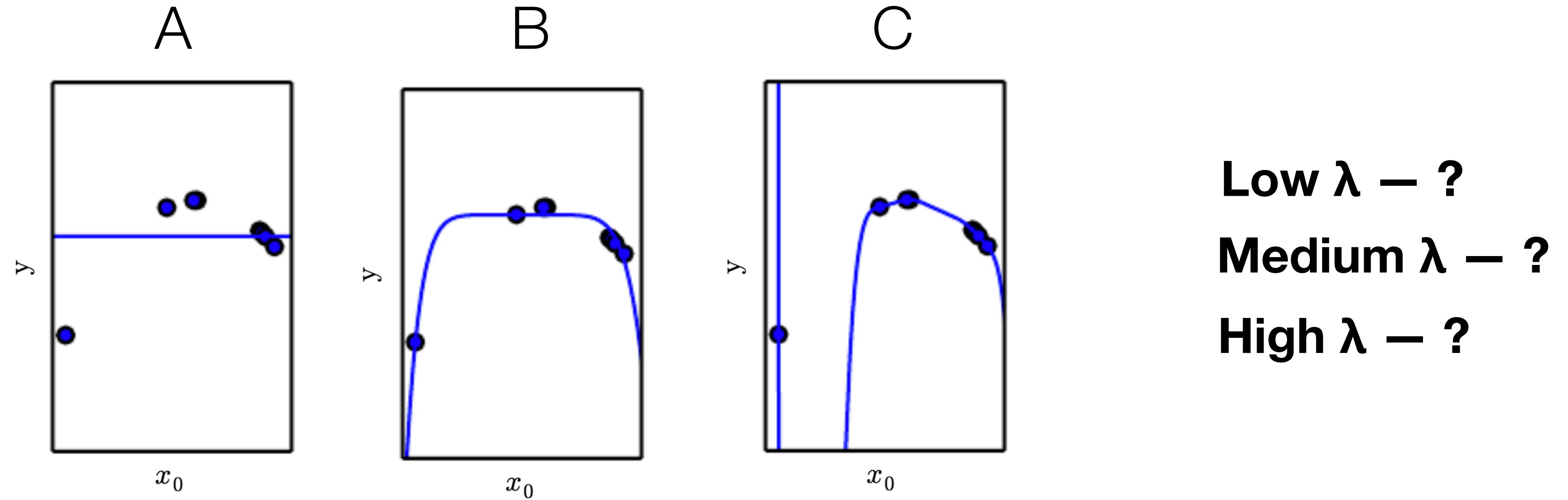
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

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

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

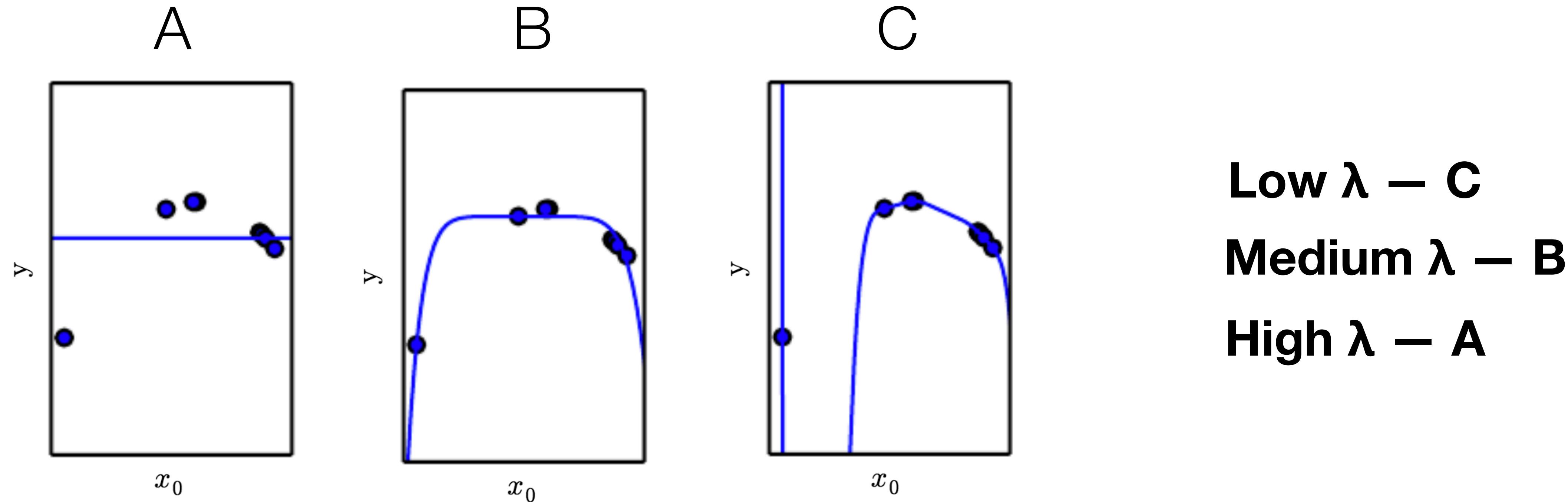
(Probabilistic interpretation: R is a Gaussian **prior** over values of the parameters.)

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



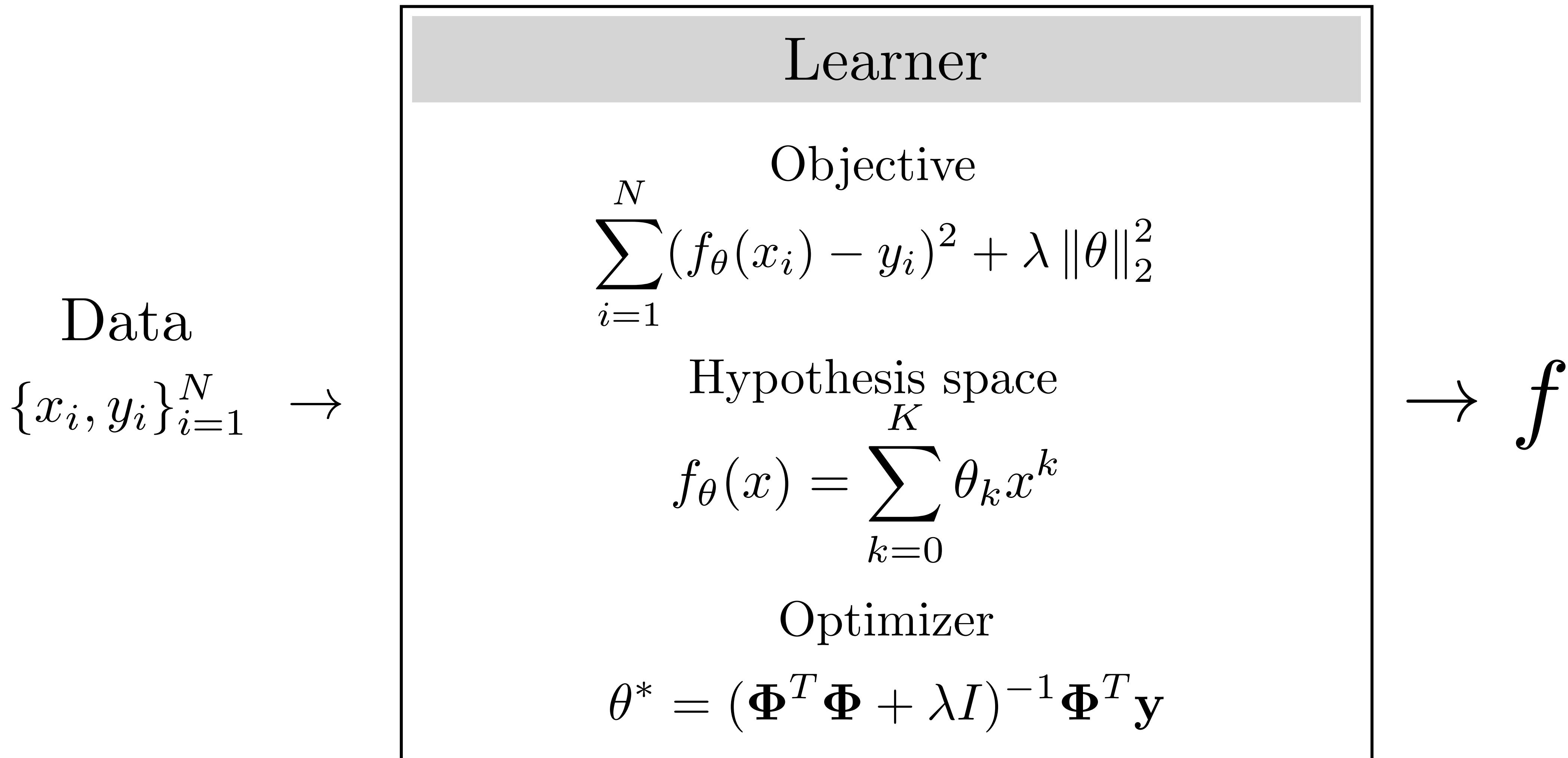
[Adapted from “Deep Learning”, Goodfellow et al.]

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



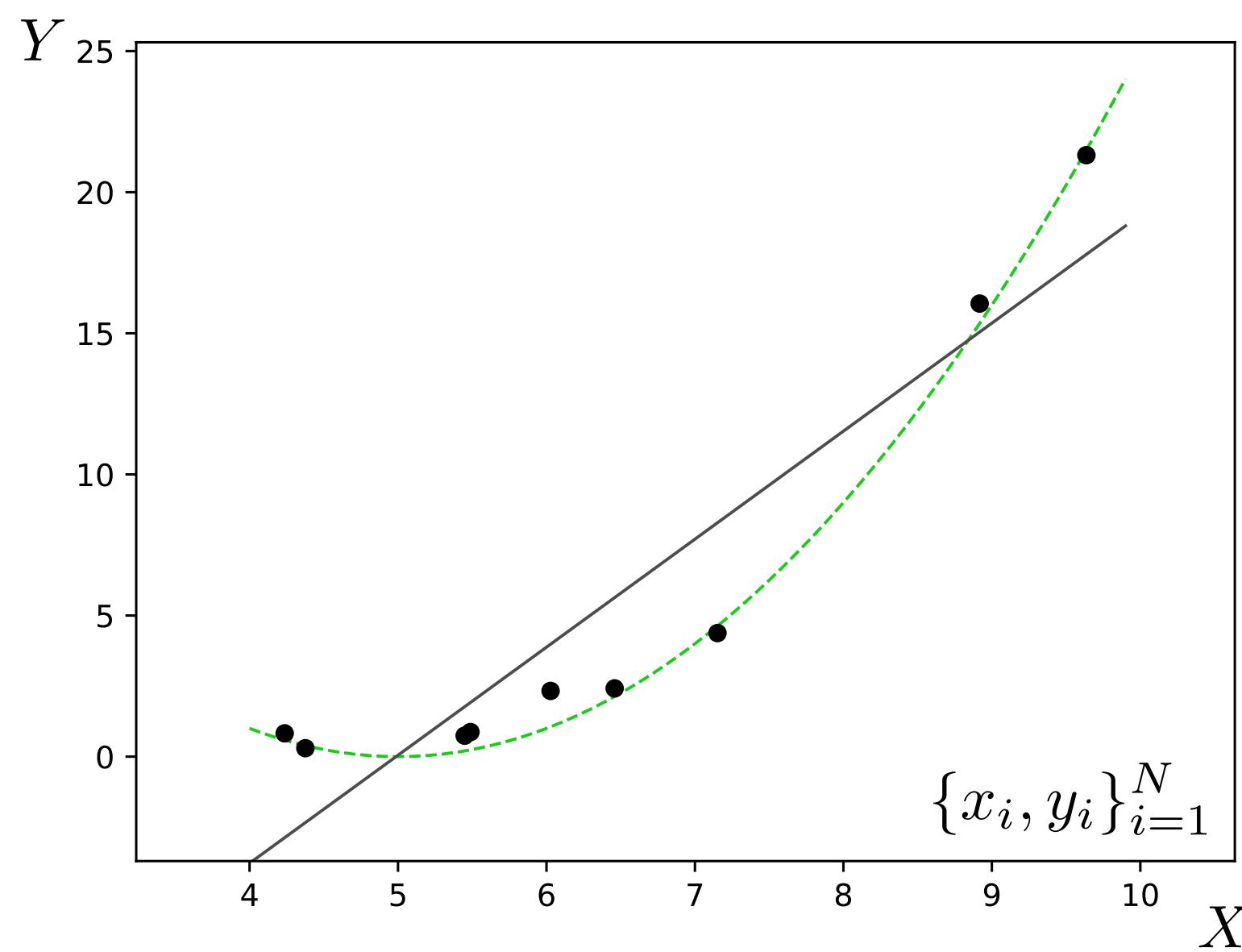
[Adapted from “Deep Learning”, Goodfellow et al.]

# Regularized polynomial least squares regression



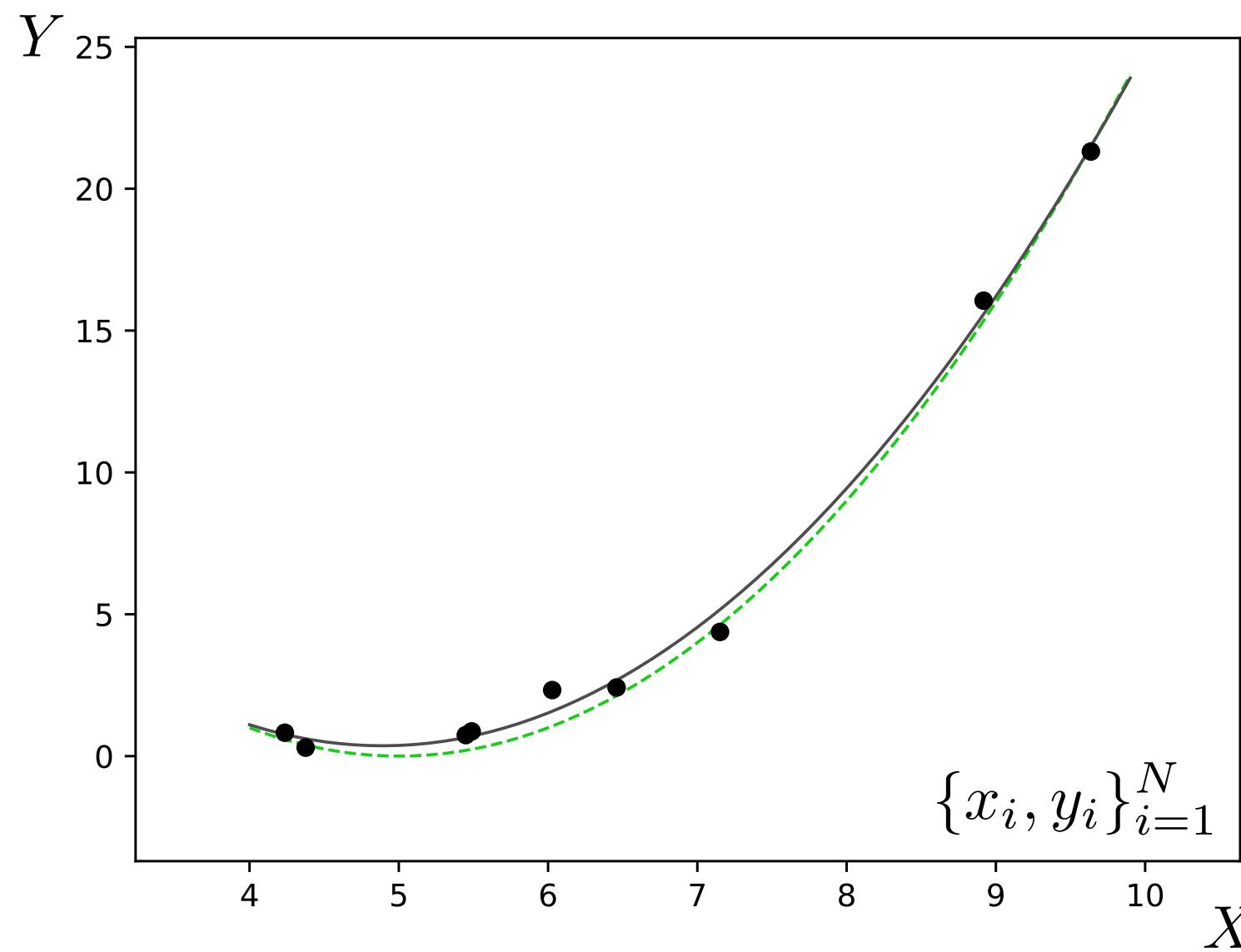
## Underfitting

$$K = 1$$



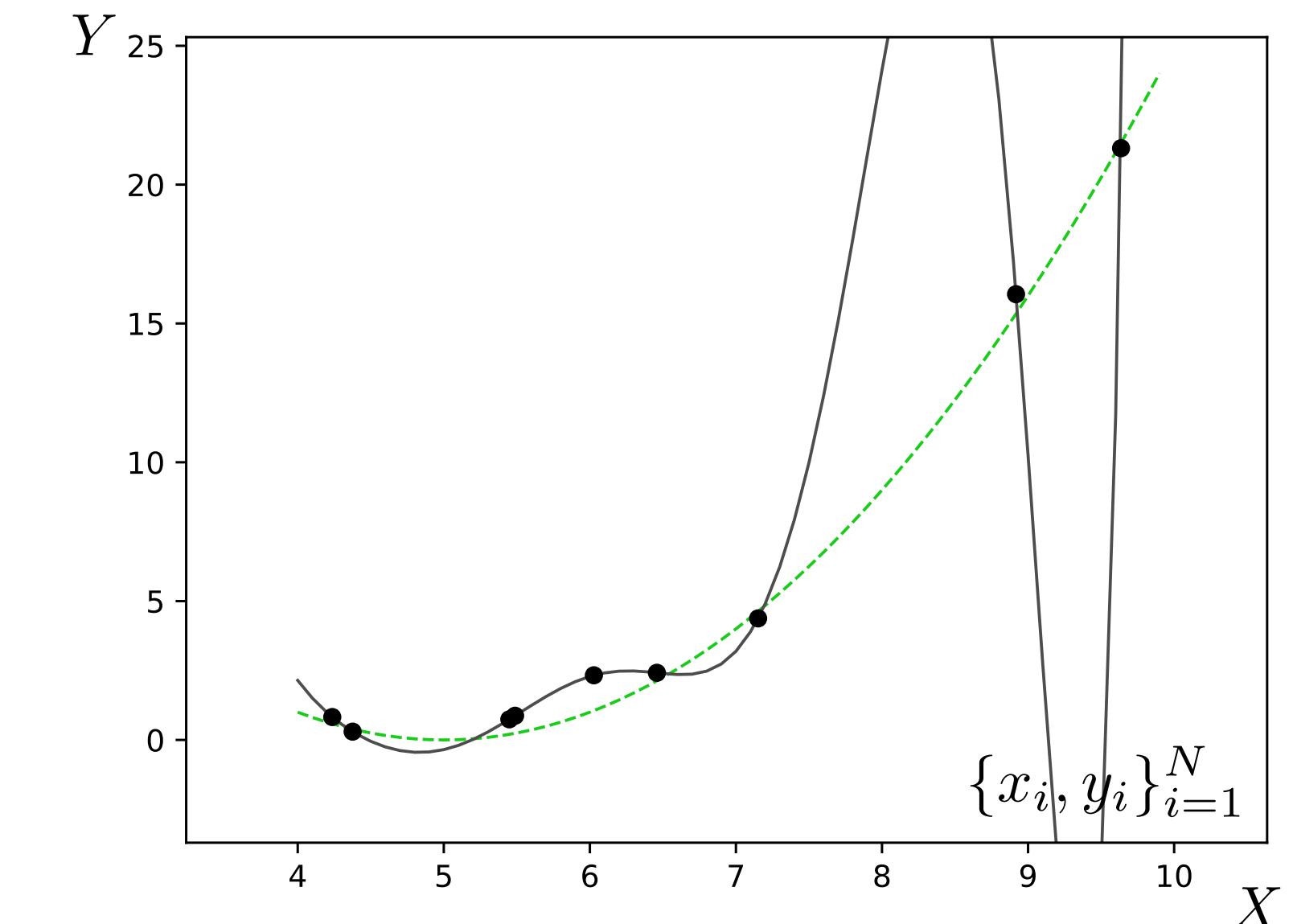
## Appropriate model

$$K = 2$$



## Overfitting

$$K = 10$$



Simple model

Doesn't fit the training data

Simple model

Fits the training data

Complex model

Fits the training data

Next lecture: neural networks