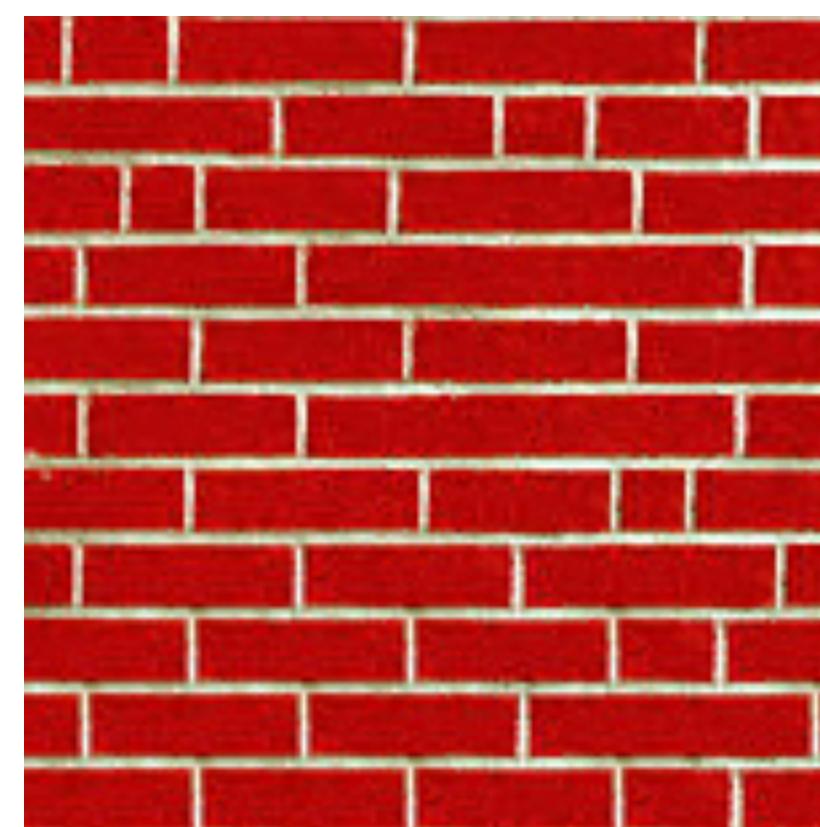


Machine learning in computer vision

Texture synthesis





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





Simple inpainting [Efros & Leung 1999]

How would a human solve this?

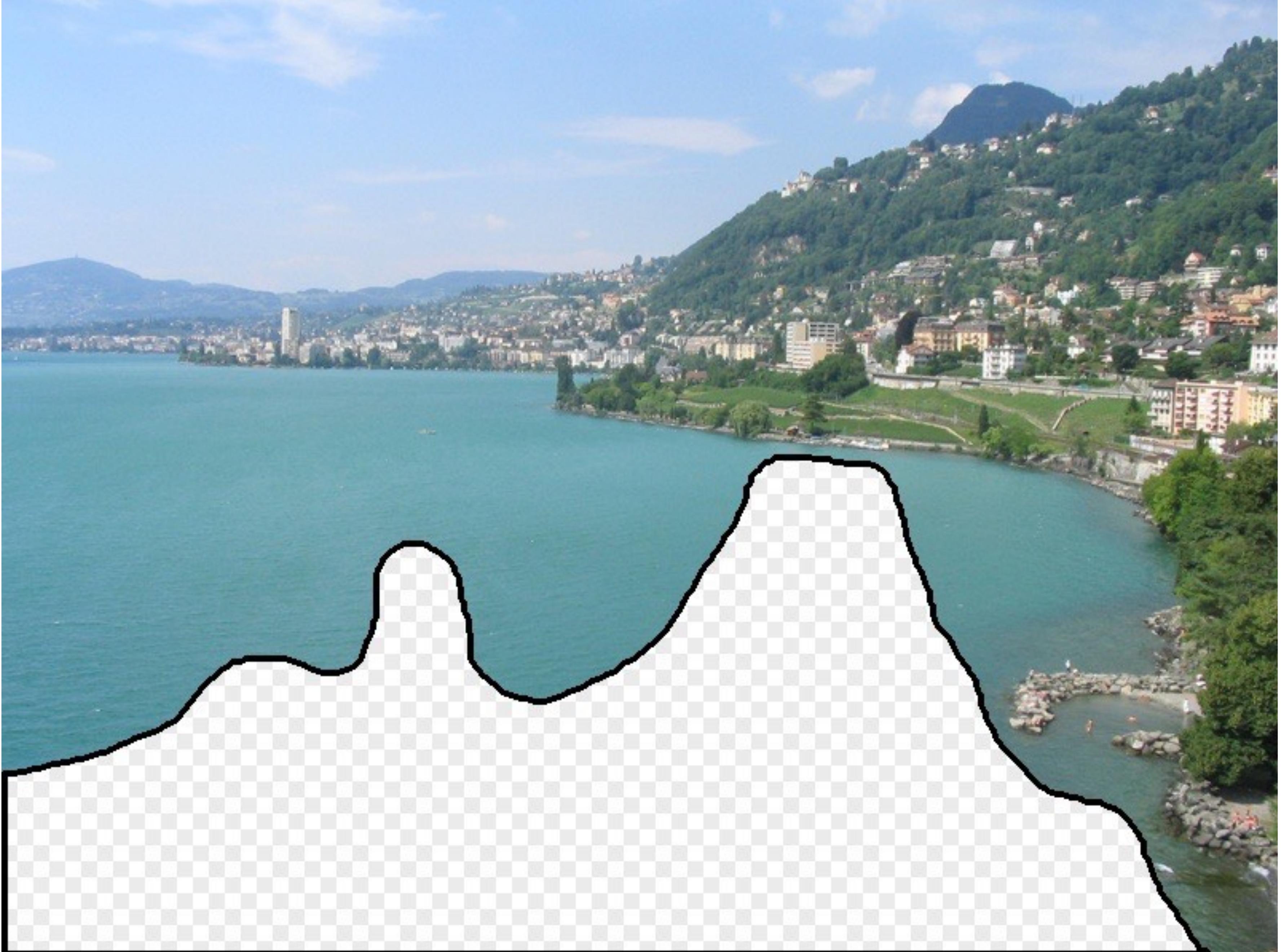




[Hays and Efros. Scene Completion Using Millions of Photographs. SIGGRAPH 2007.]

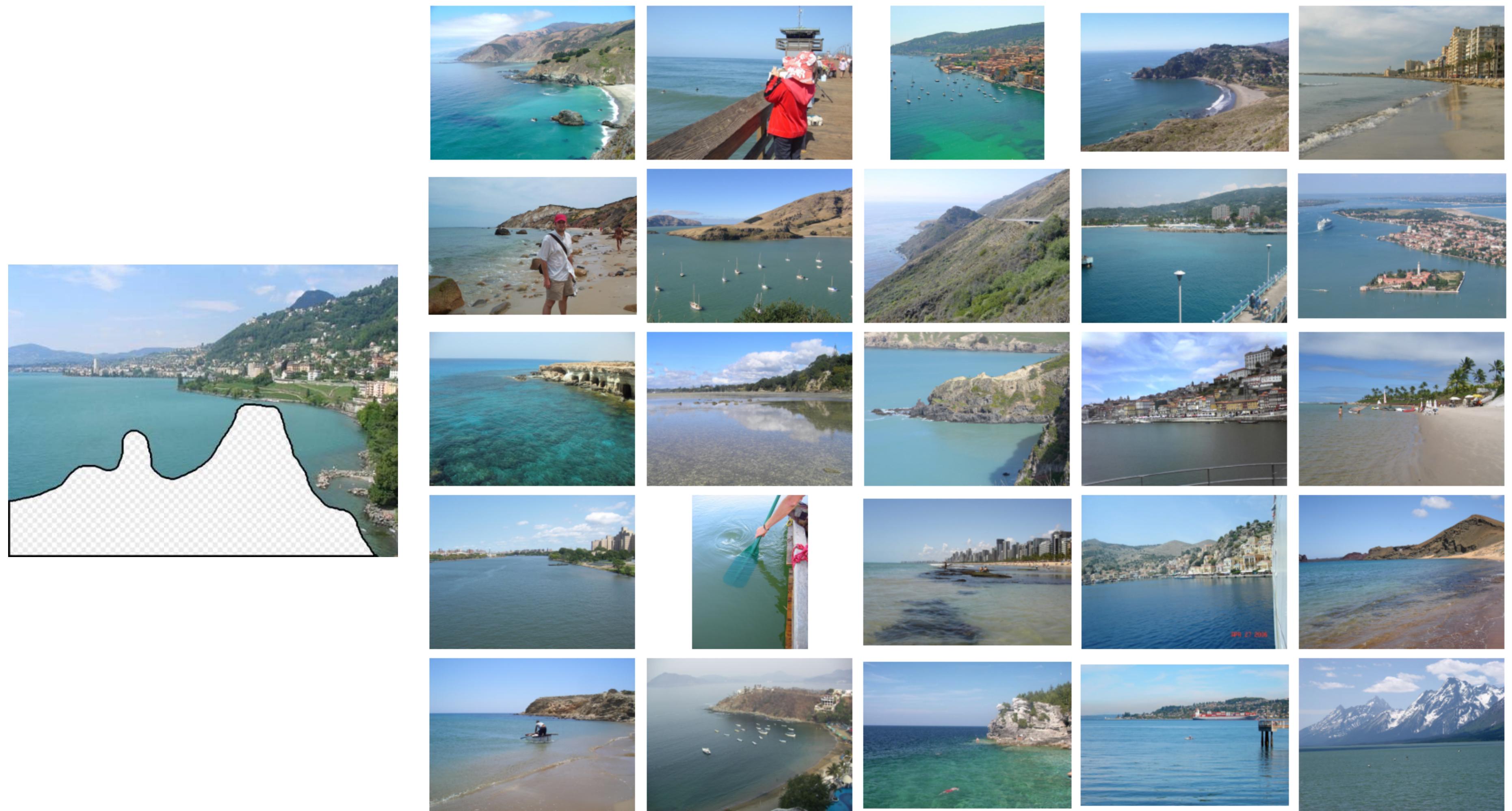
7

Source: A. Efros

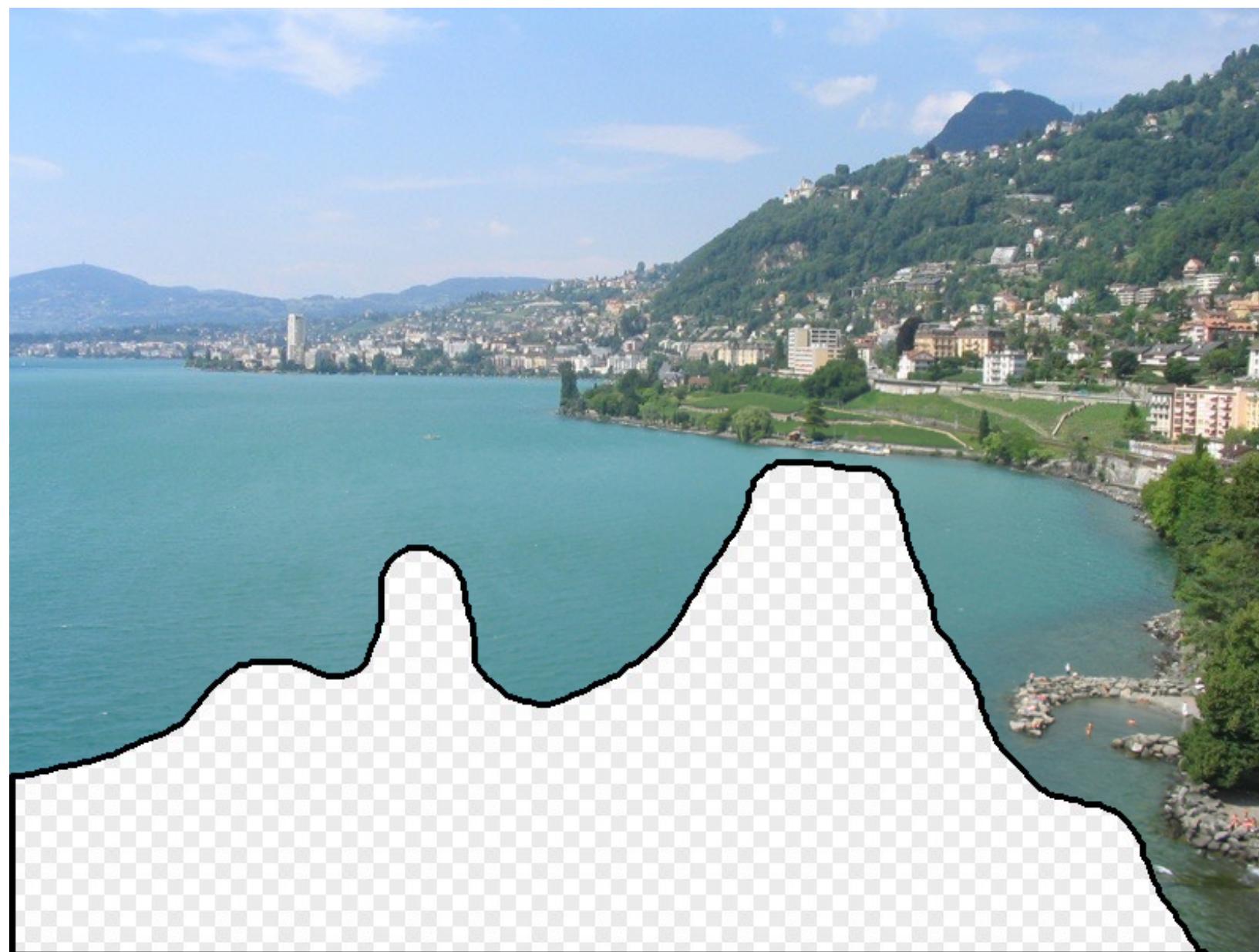


2 Million Flickr Images





... 200 total





12

Source: A. Efros

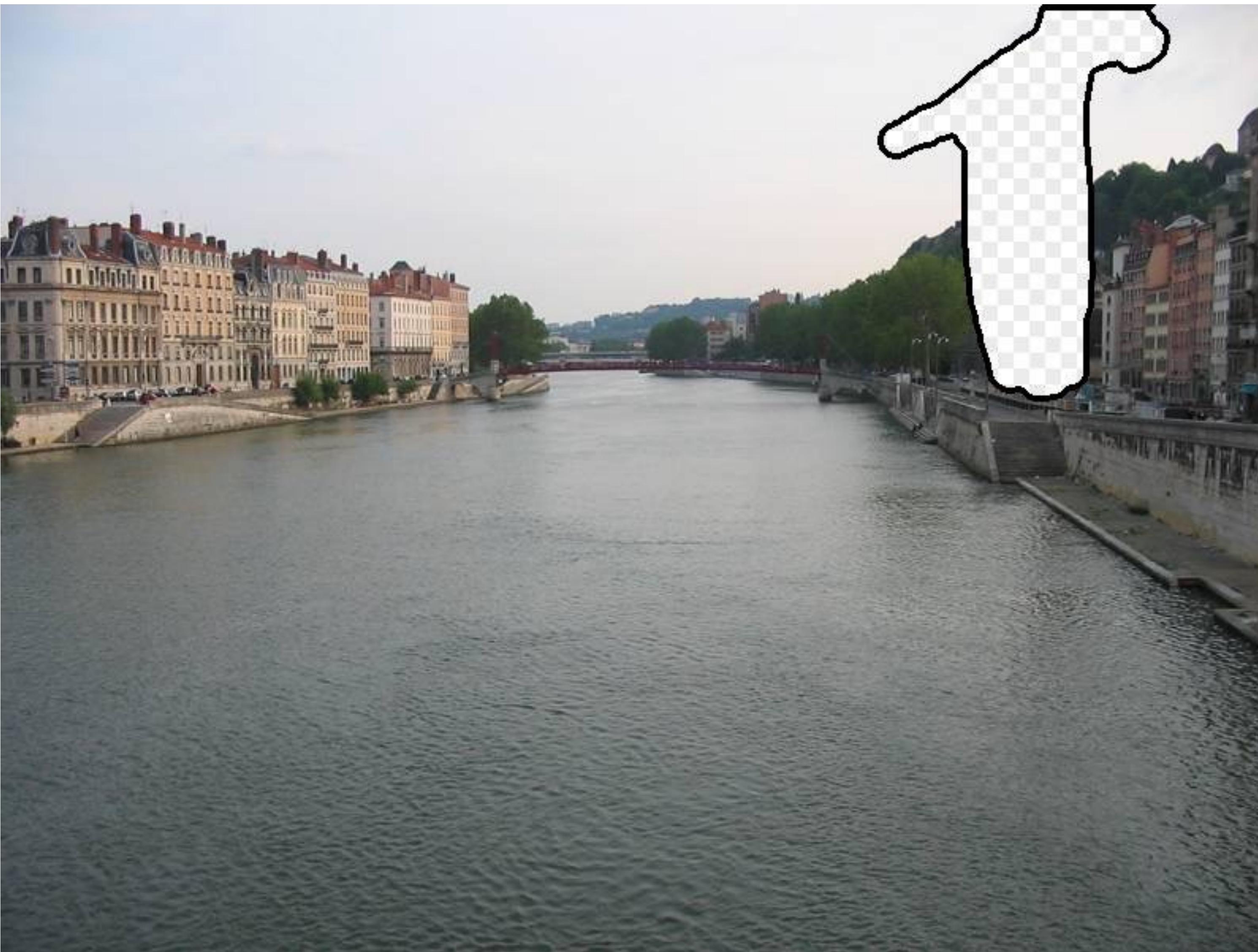




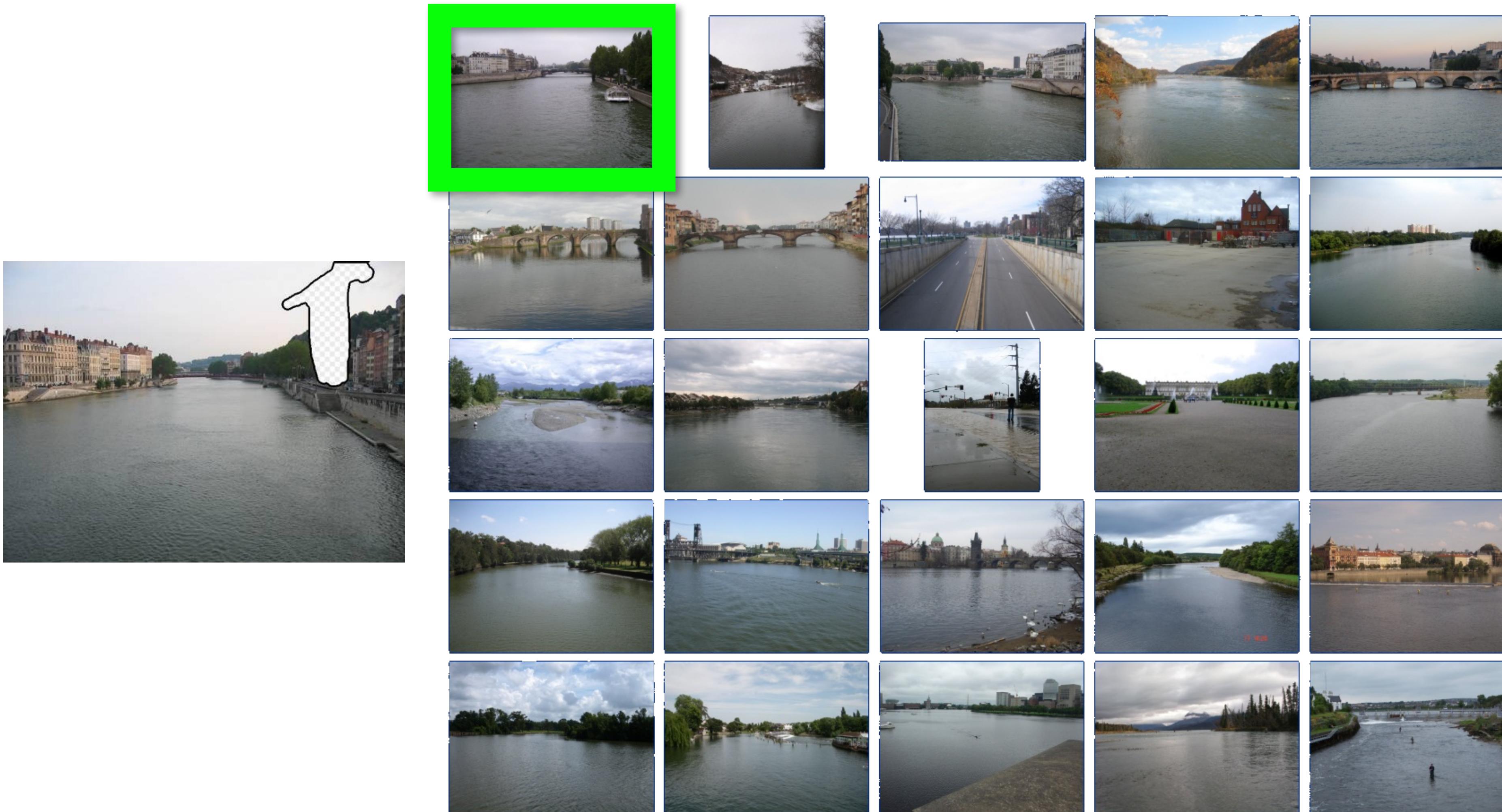












... 200 scene matches

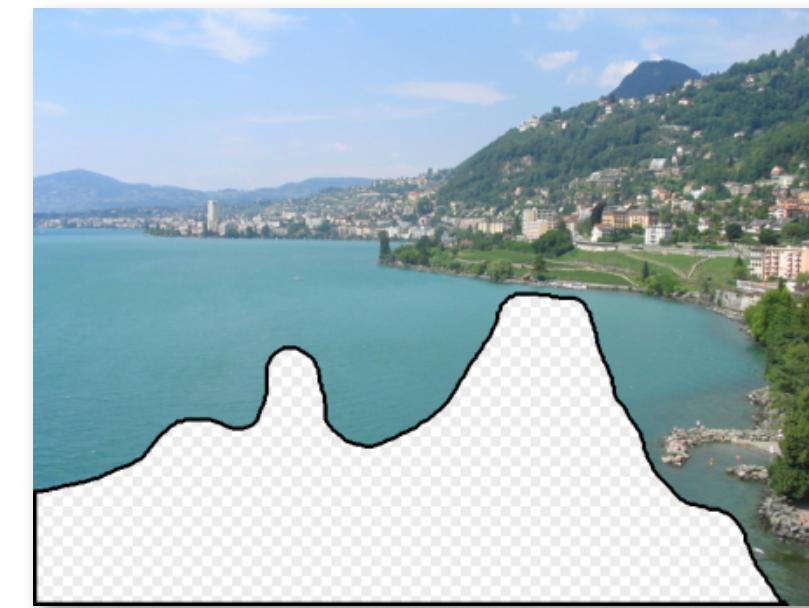


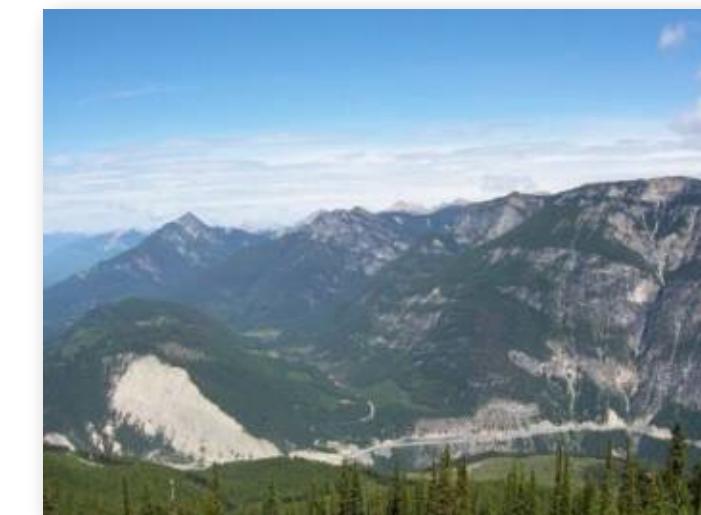
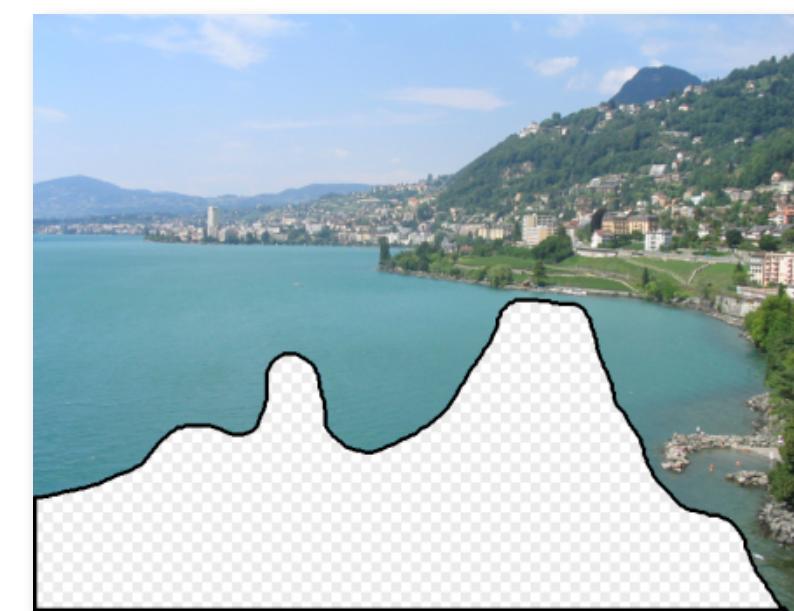




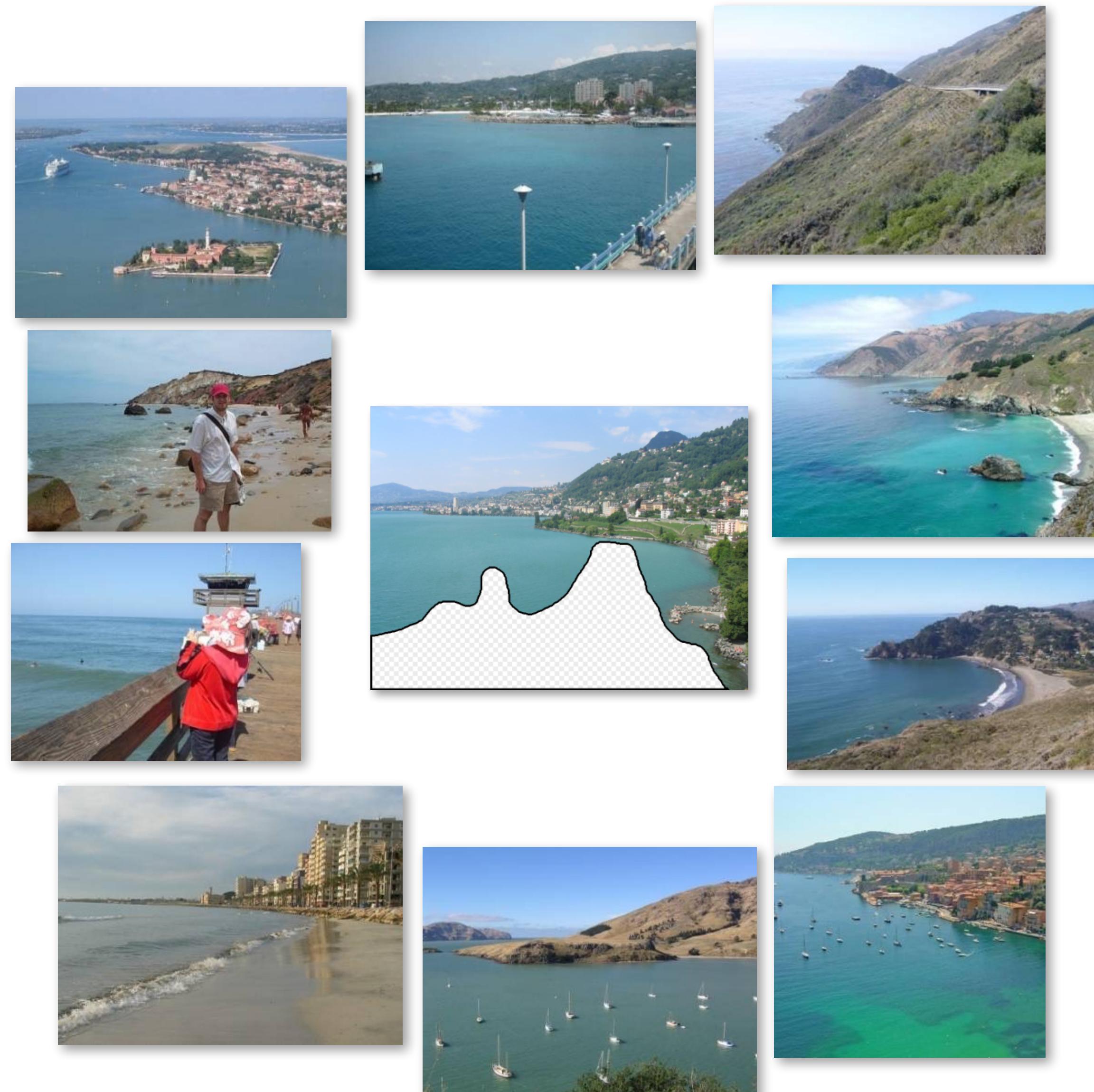


Why does this work?





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]

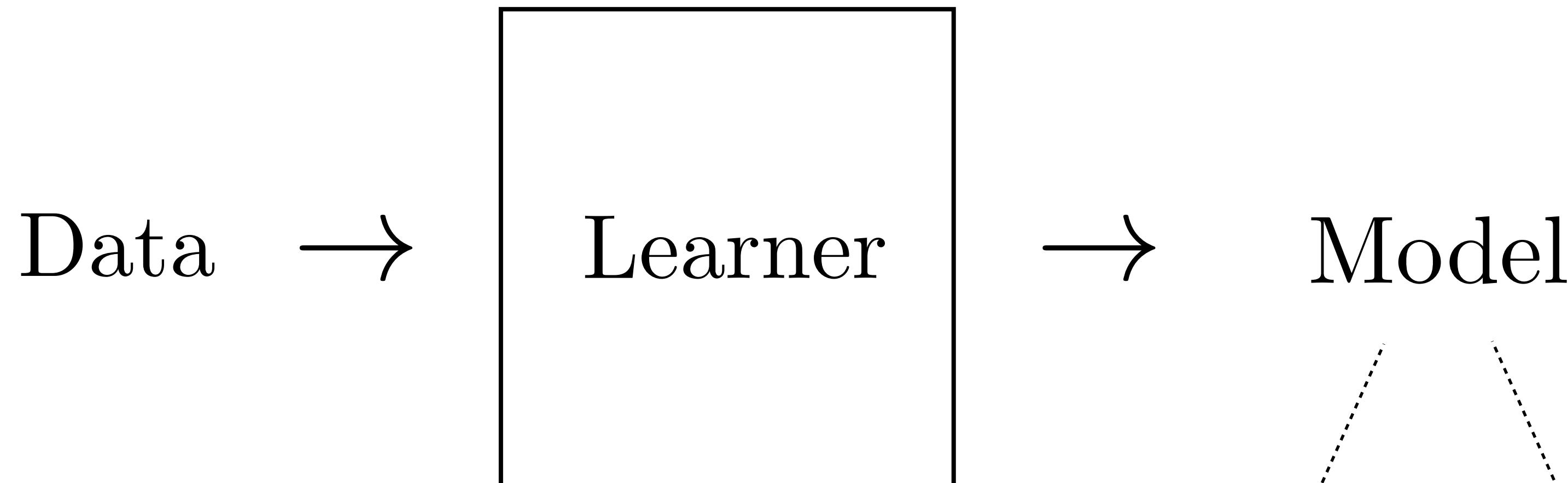
Parts of our world can be explained by elegant
mathematics
physics, chemistry, astronomy, etc.

But much cannot
psychology, economics, genetics, etc.

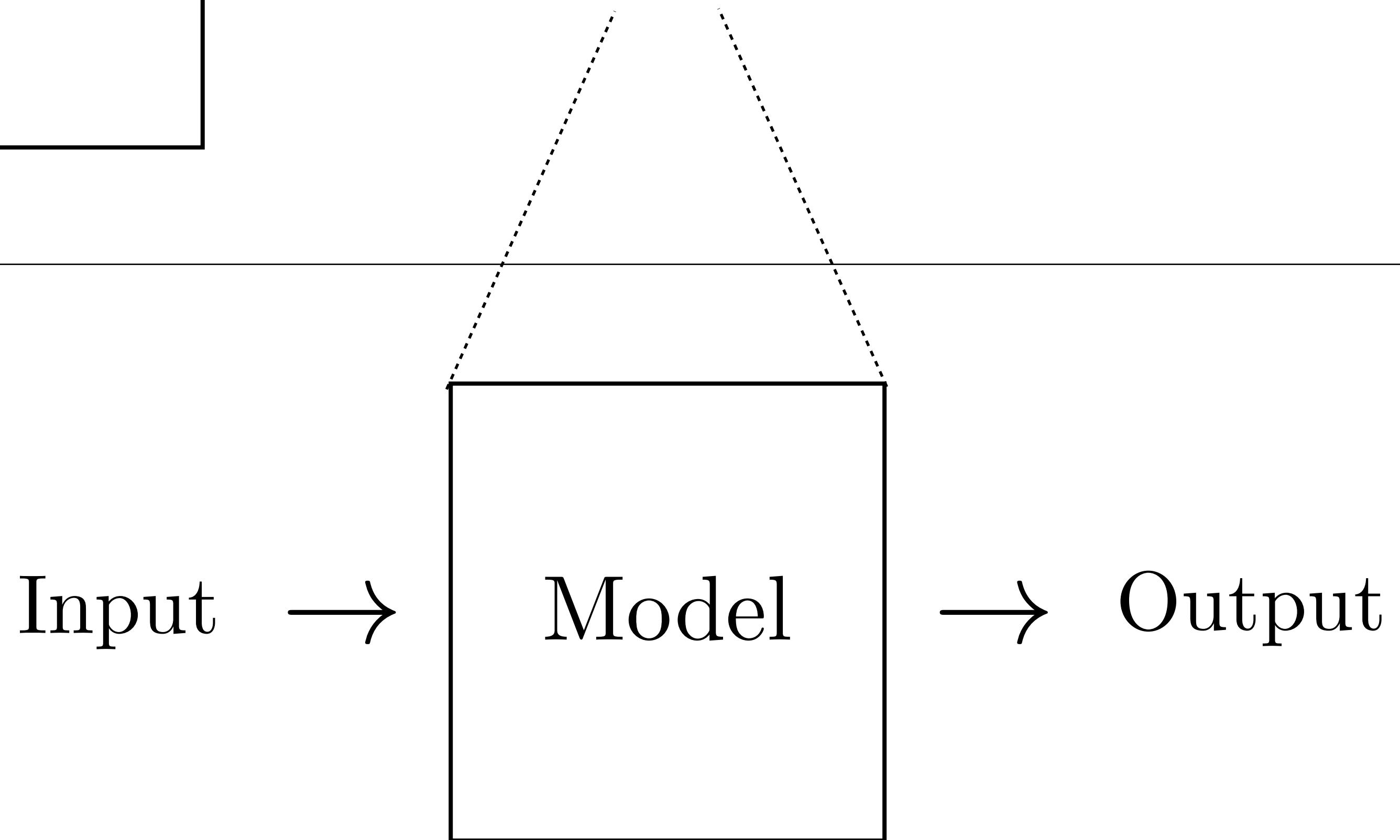
Enter The Data!

“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**)

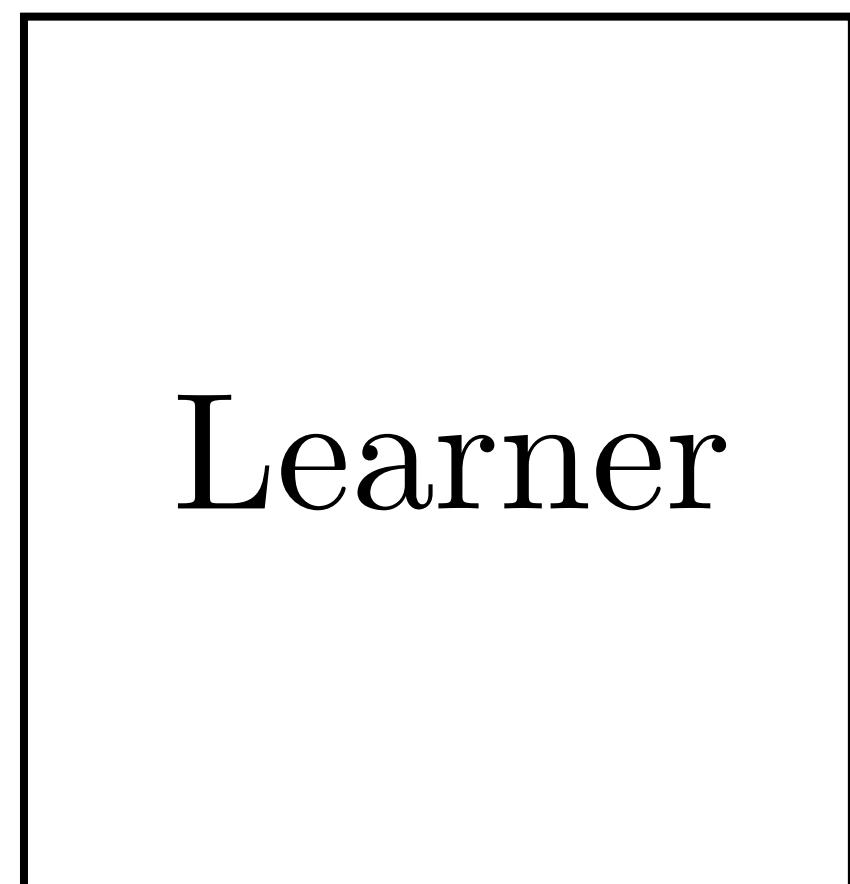
Training data

$$\{x_1, y_1\}$$

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

$$\{x_3, y_3\}$$

...

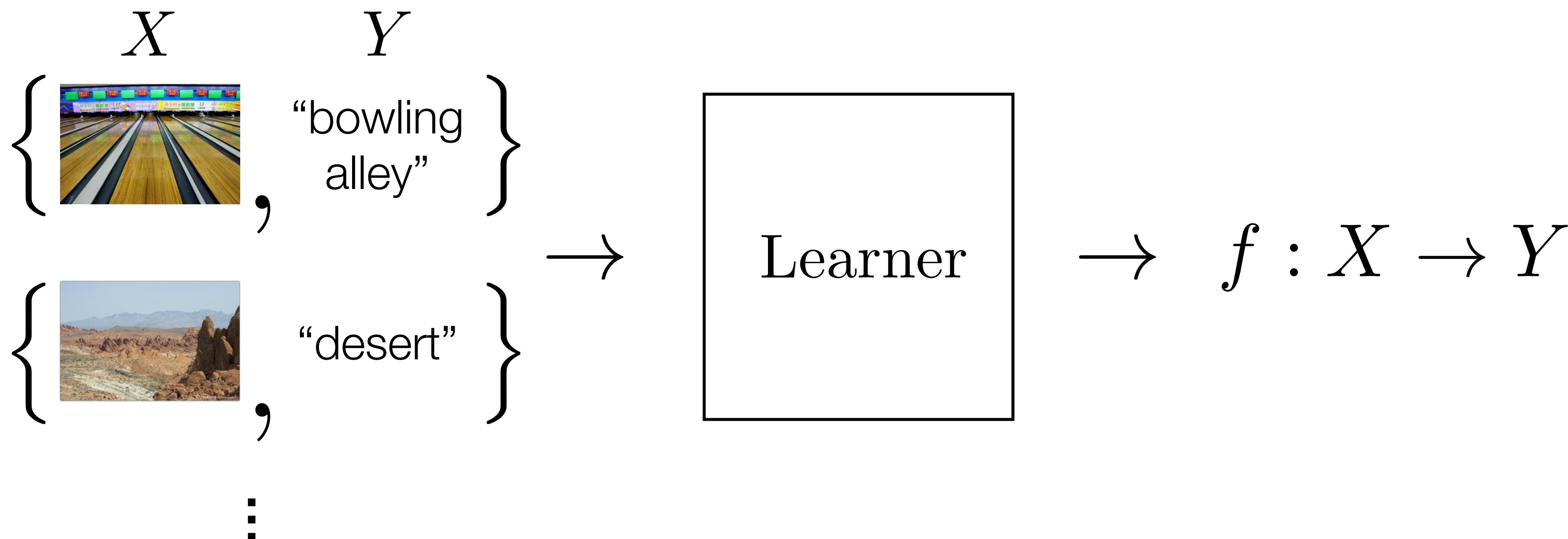


$$\rightarrow f : X \rightarrow Y$$

Learning from examples

(aka **supervised learning**)

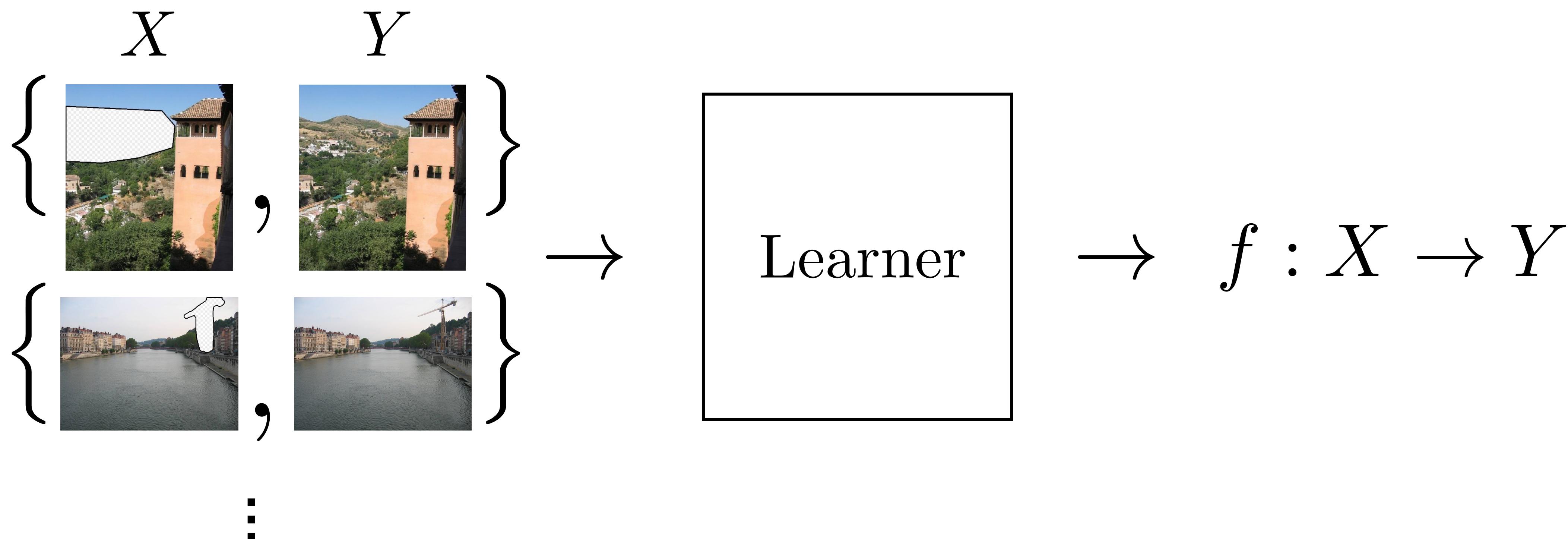
Training data



Learning from examples

(aka **supervised learning**)

Training data



Case study #1: Nearest neighbor

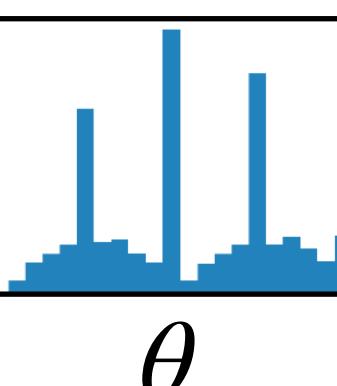
Nearest neighbor



Input image

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

Feature vector x_q



Some feature ideas:

tiny image

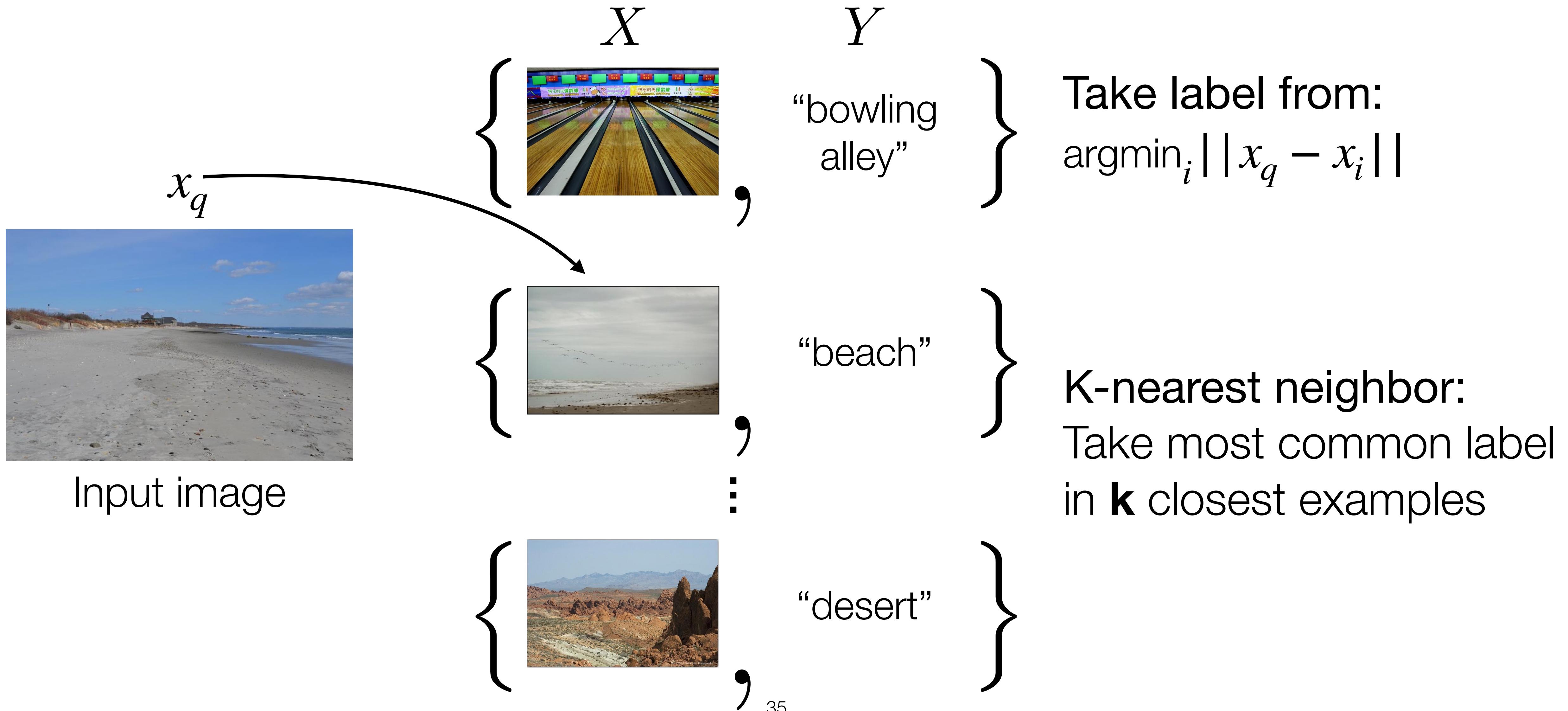
normalized
tiny image

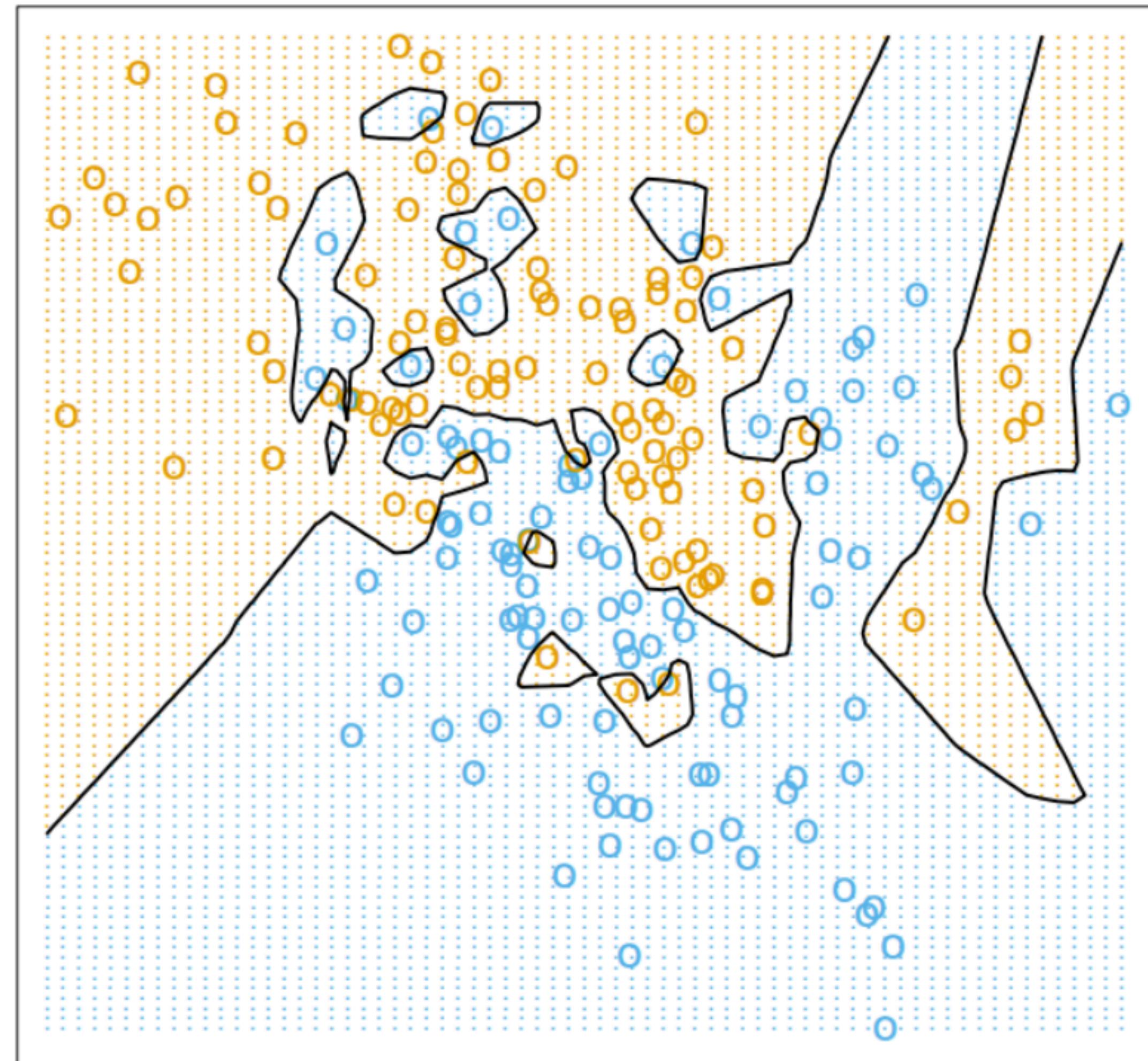
$$\frac{x - \mu}{\sigma}$$

edge orientation
histogram

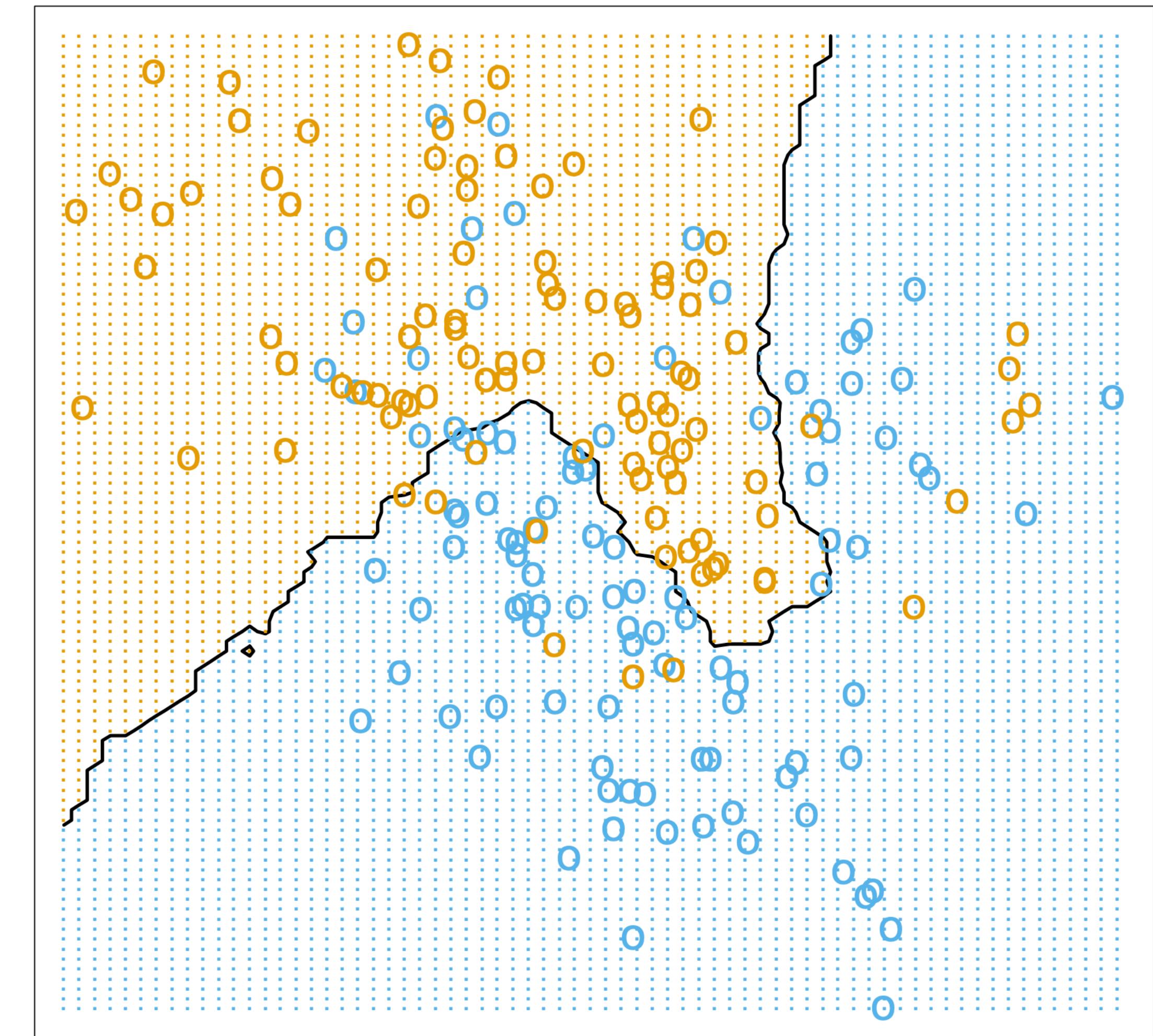
We'll discuss better options
later in course.

Nearest neighbor





$k = 1$



$k = 15$

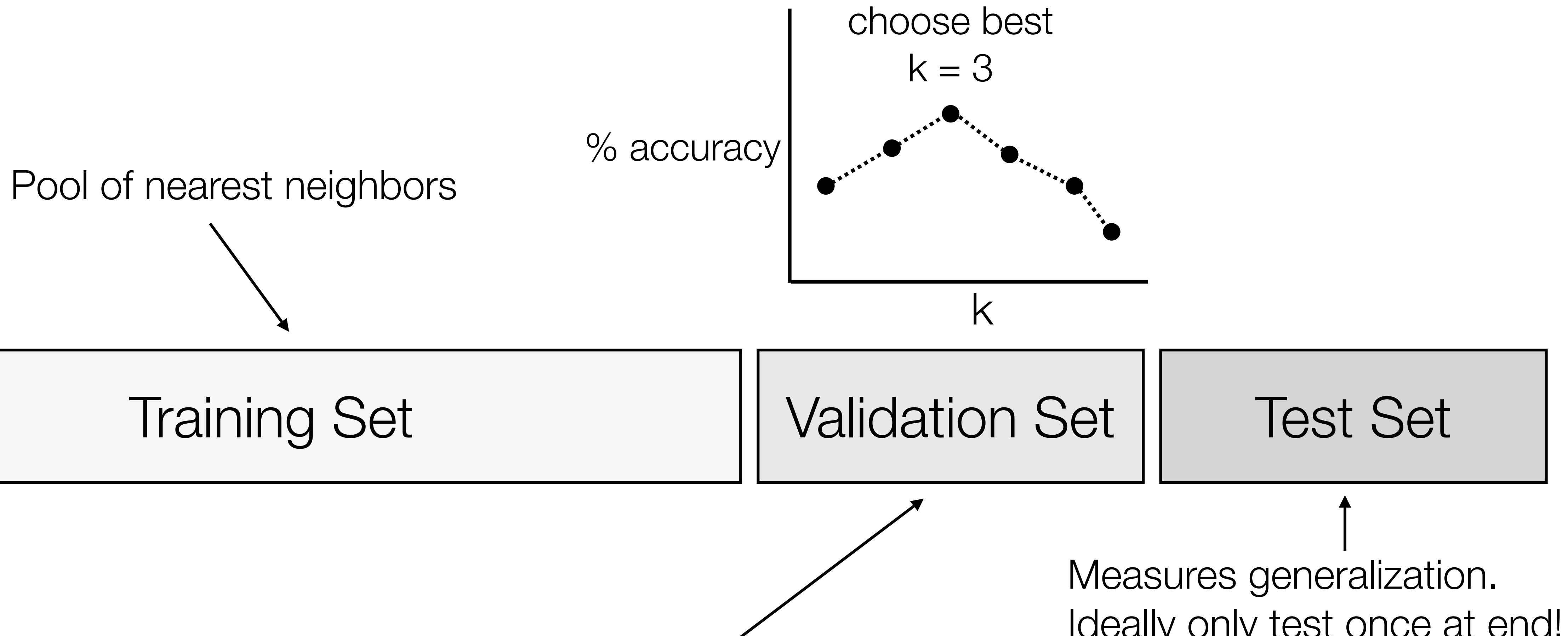
Source: [Hastie et al. "Elements of Statistical Learning", 2001]

How do we choose k?

Training Data

Test Data

How do we choose k ?

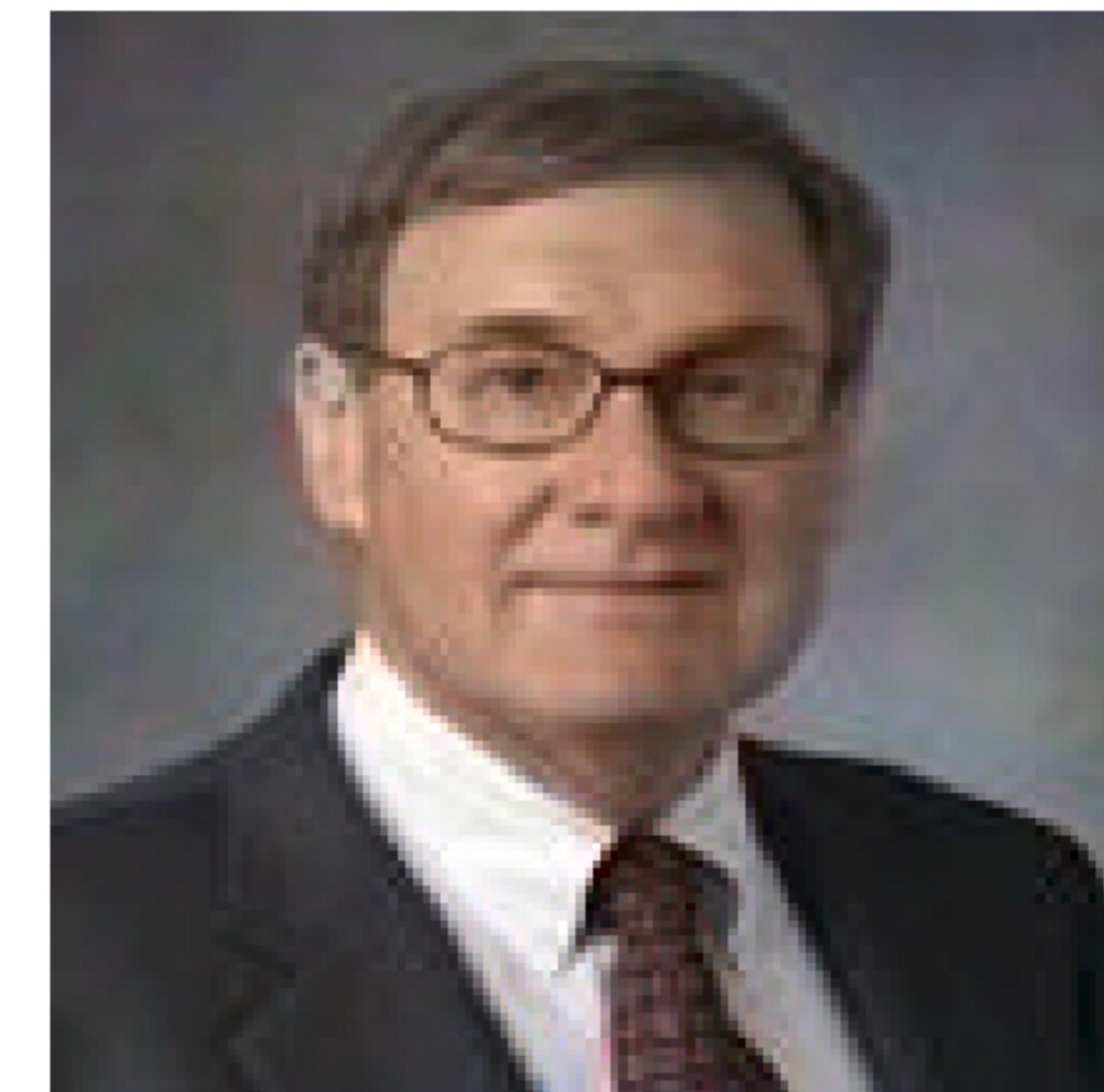


Choose **hyperparameters** like k , feature space, similarity function, etc.

Can work well when datasets are really big



Input image



Colorized output



Nearest Neighbors

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

Nearest Neighbor

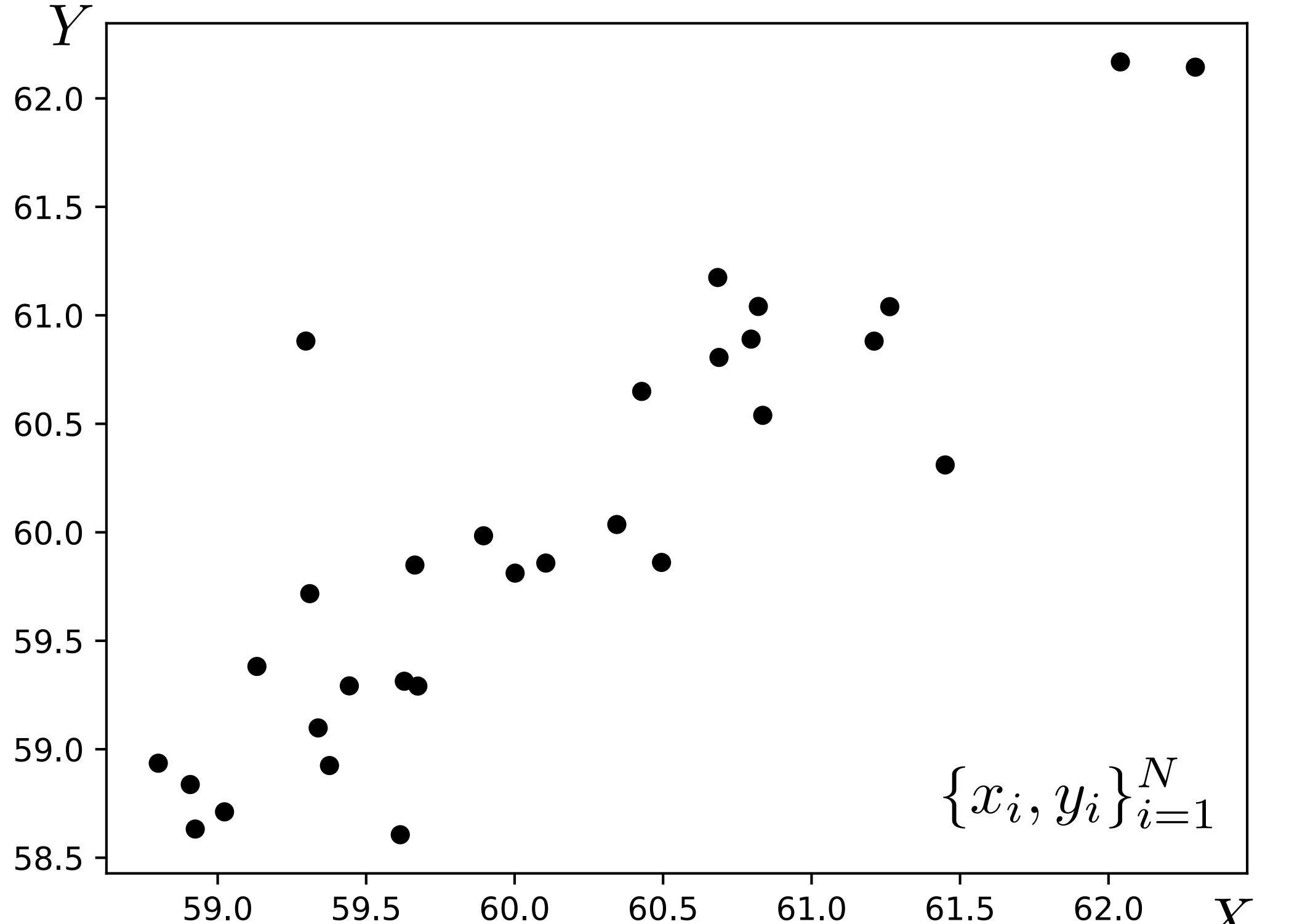
- Hard to define similarity metric



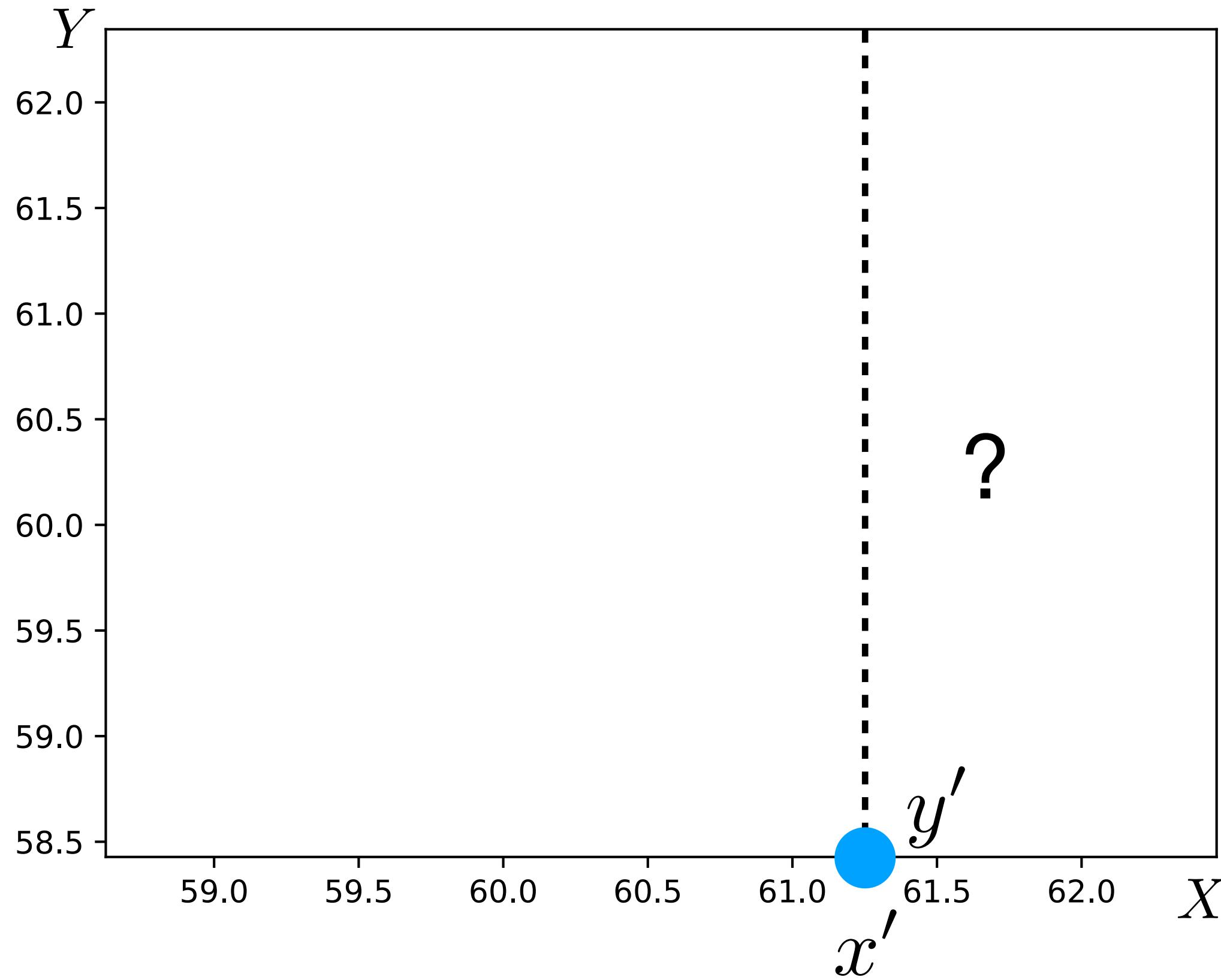
- Pays attention to all features equally
- Can be very slow: $O(nd)$

Case study #2: Linear least squares

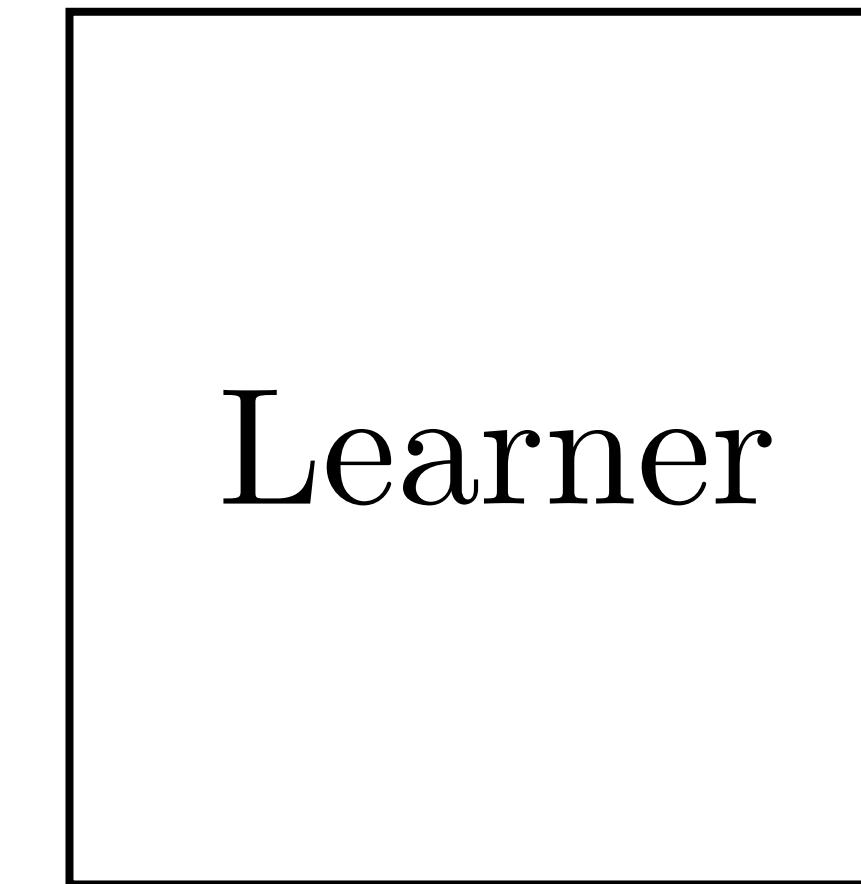
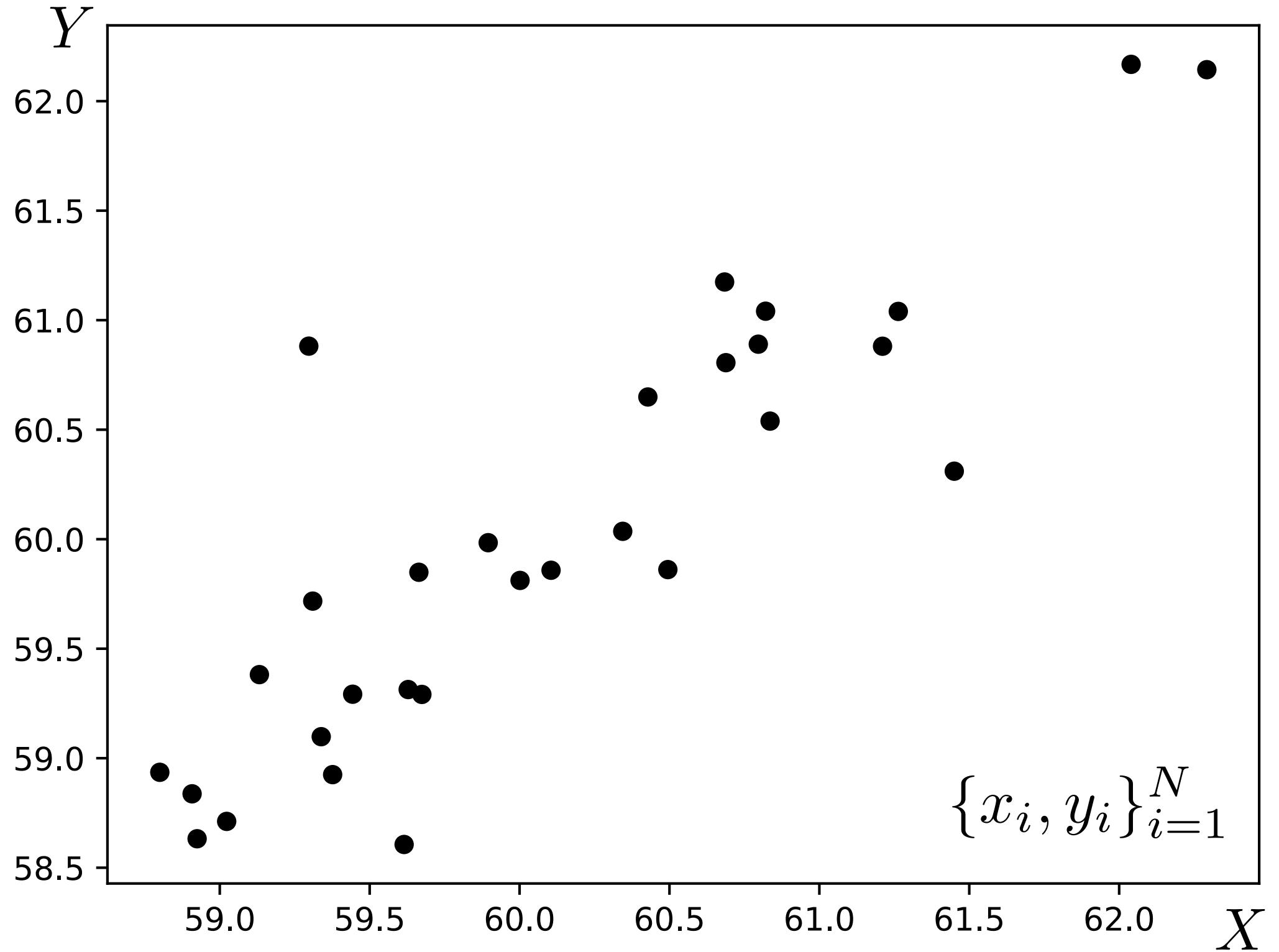
Training data



Test query



Training data

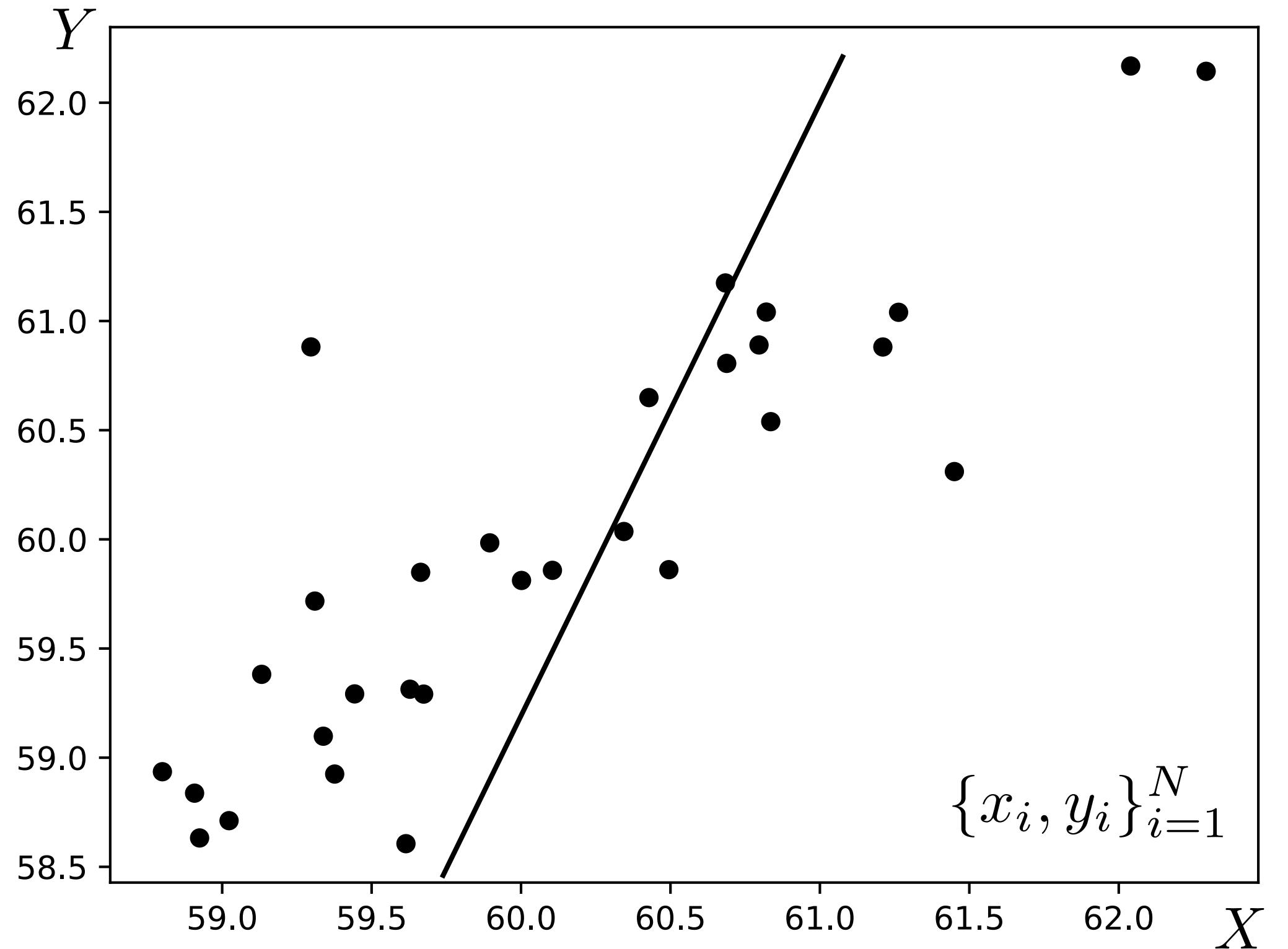


$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Hypothesis space

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

Training data

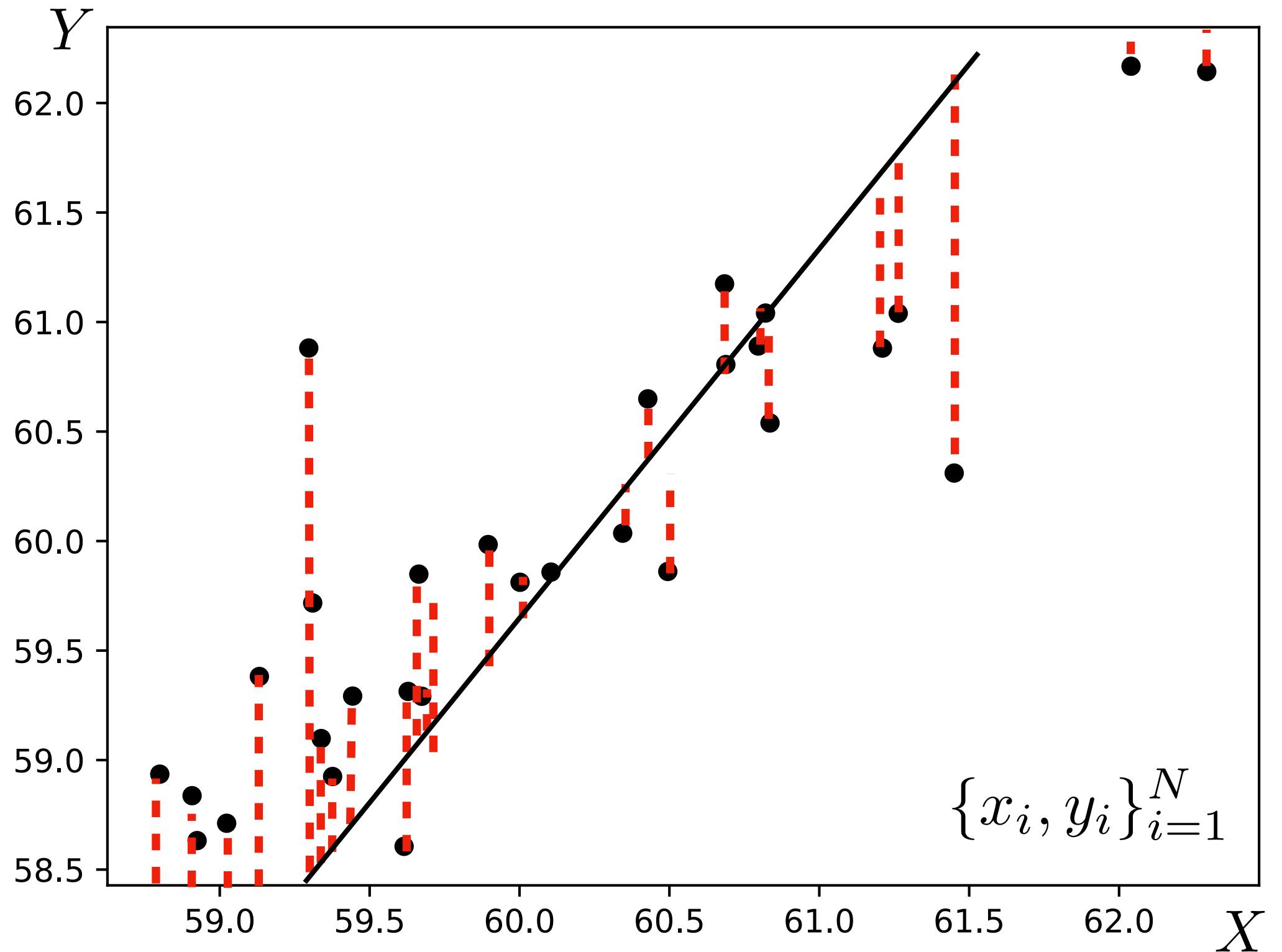


Search for the **parameters**, $\theta = \{\theta_0, \theta_1\}$, that best fit the data.

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

Best fit in what sense?

Training data



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

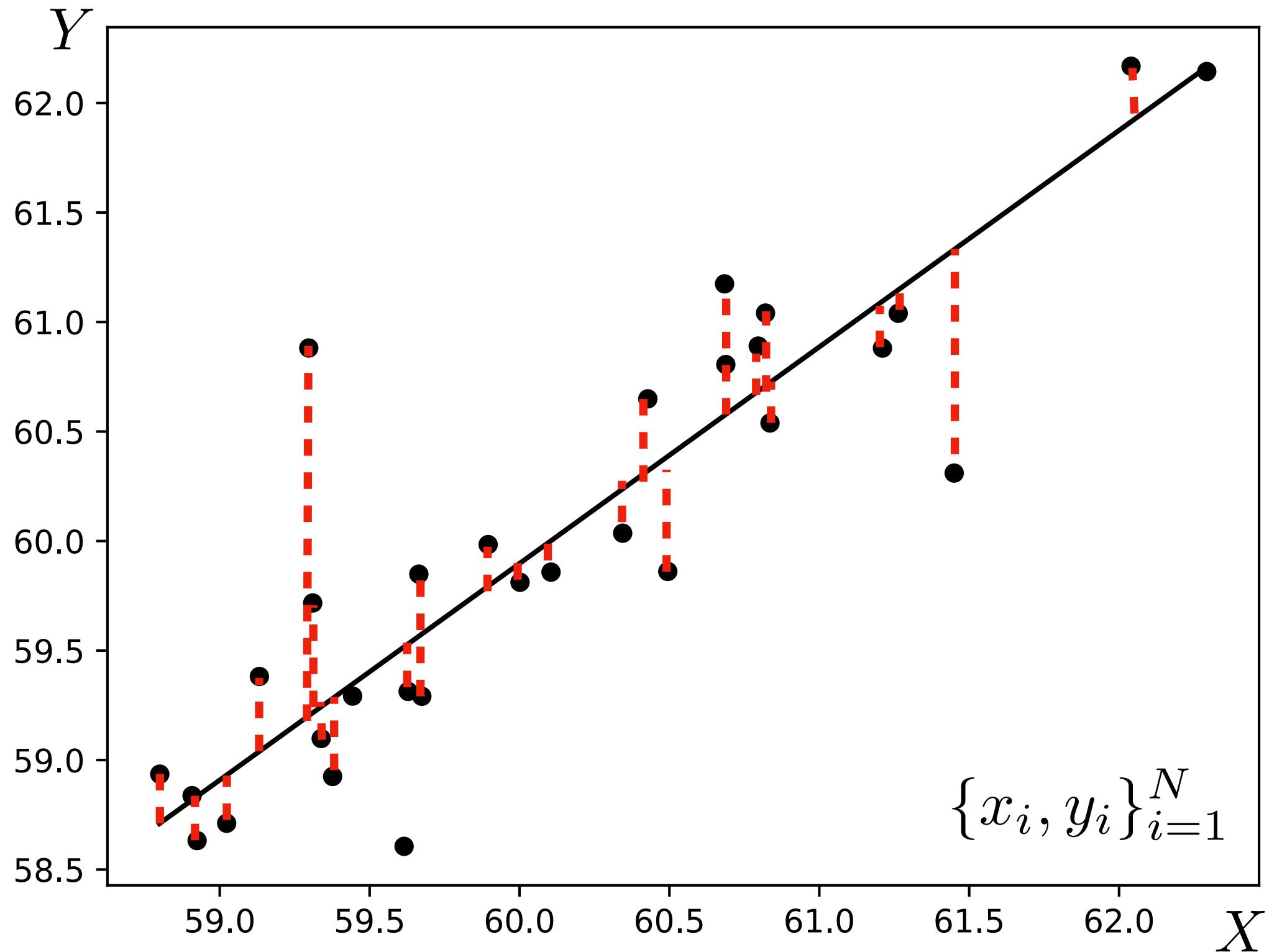
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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.

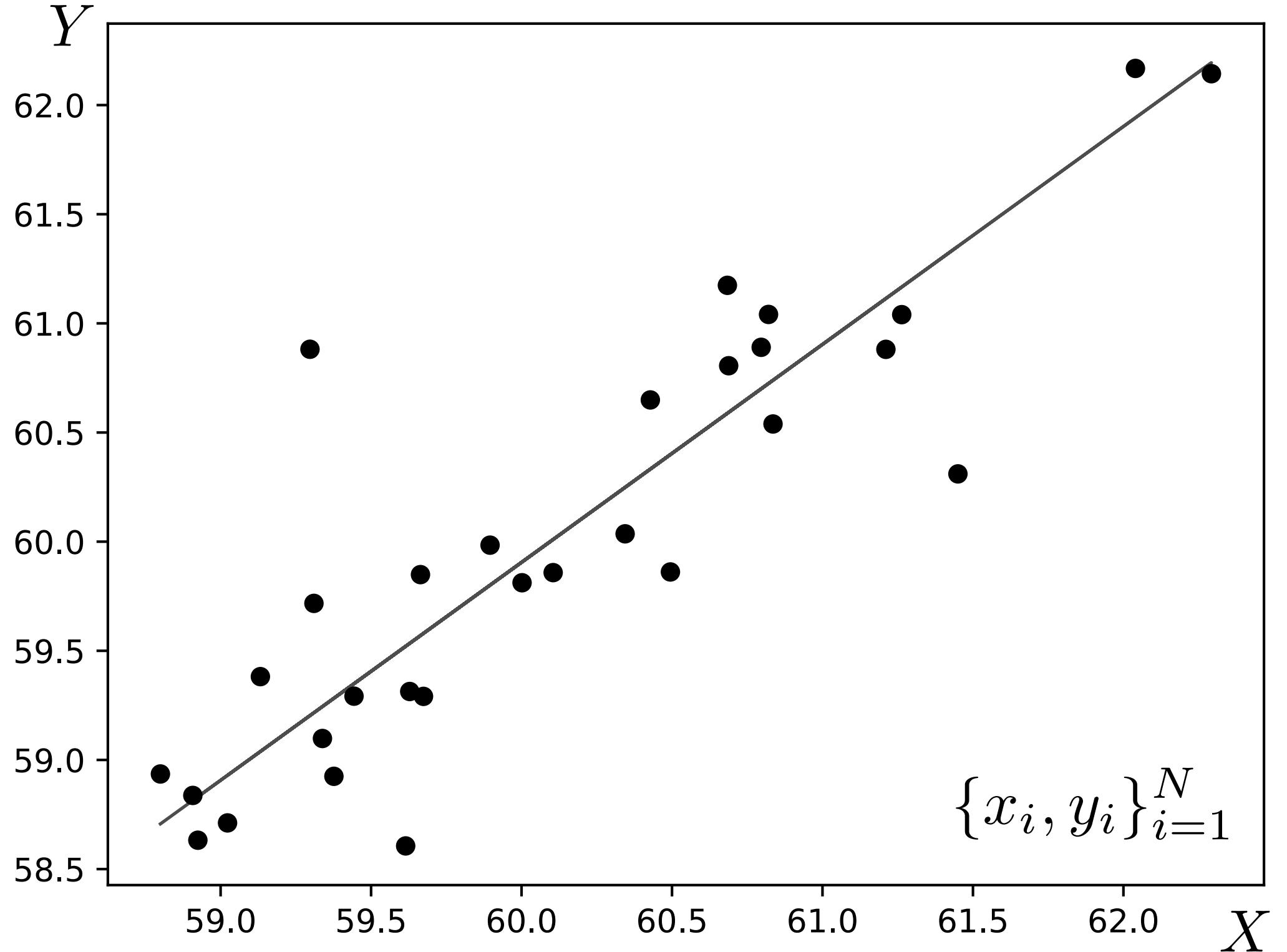
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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

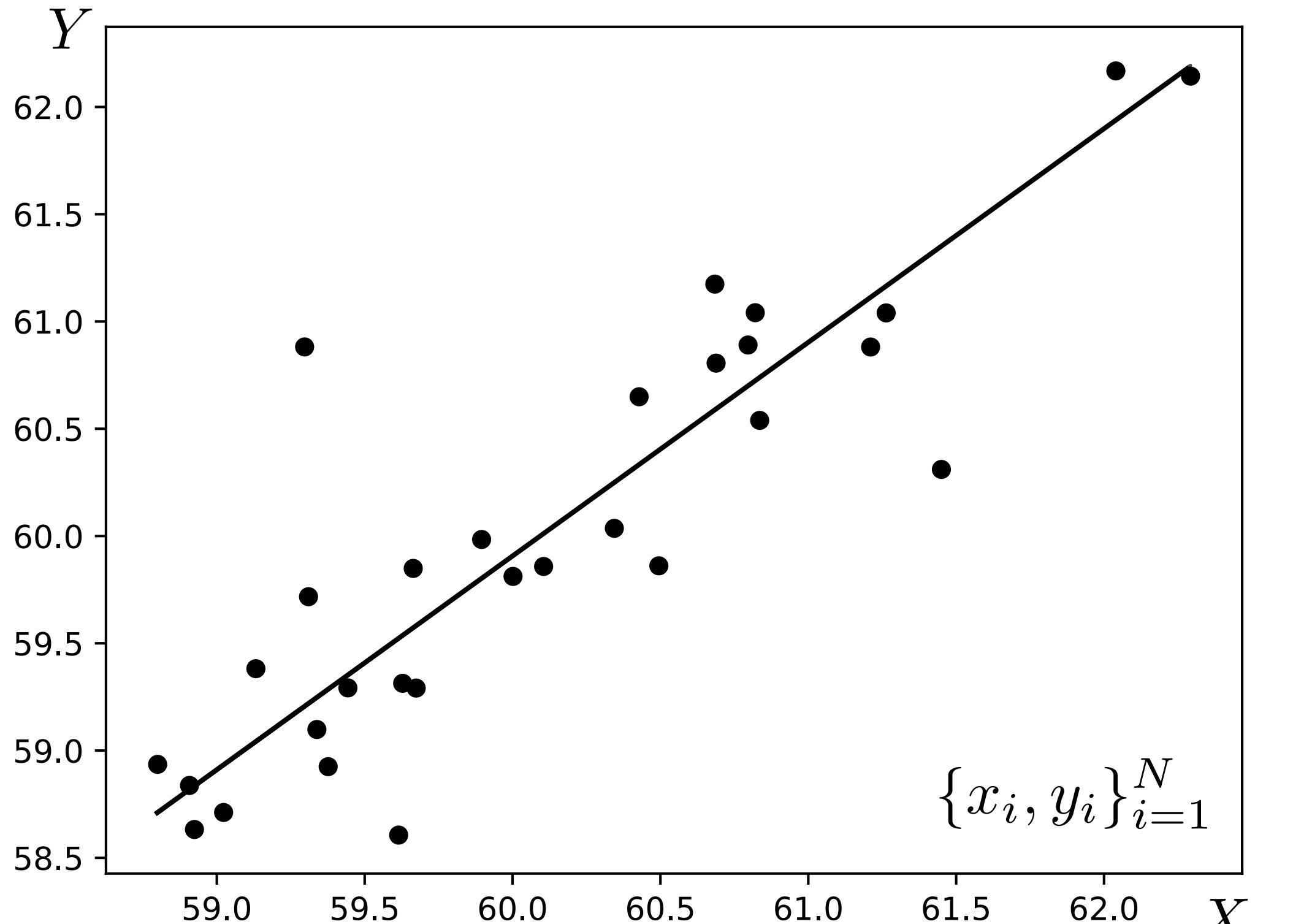


Complete learning problem:

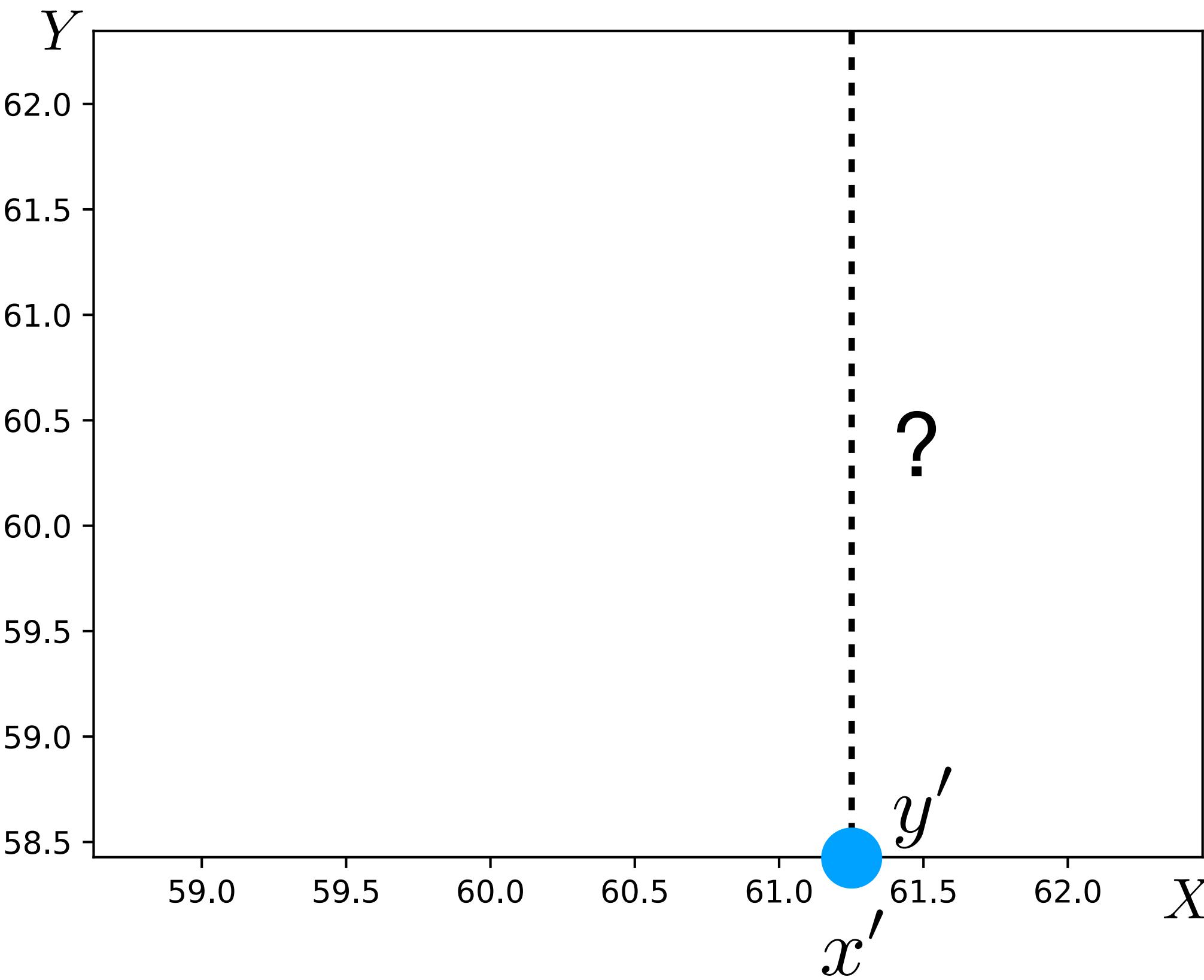
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x_i) - y_i)^2$$

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

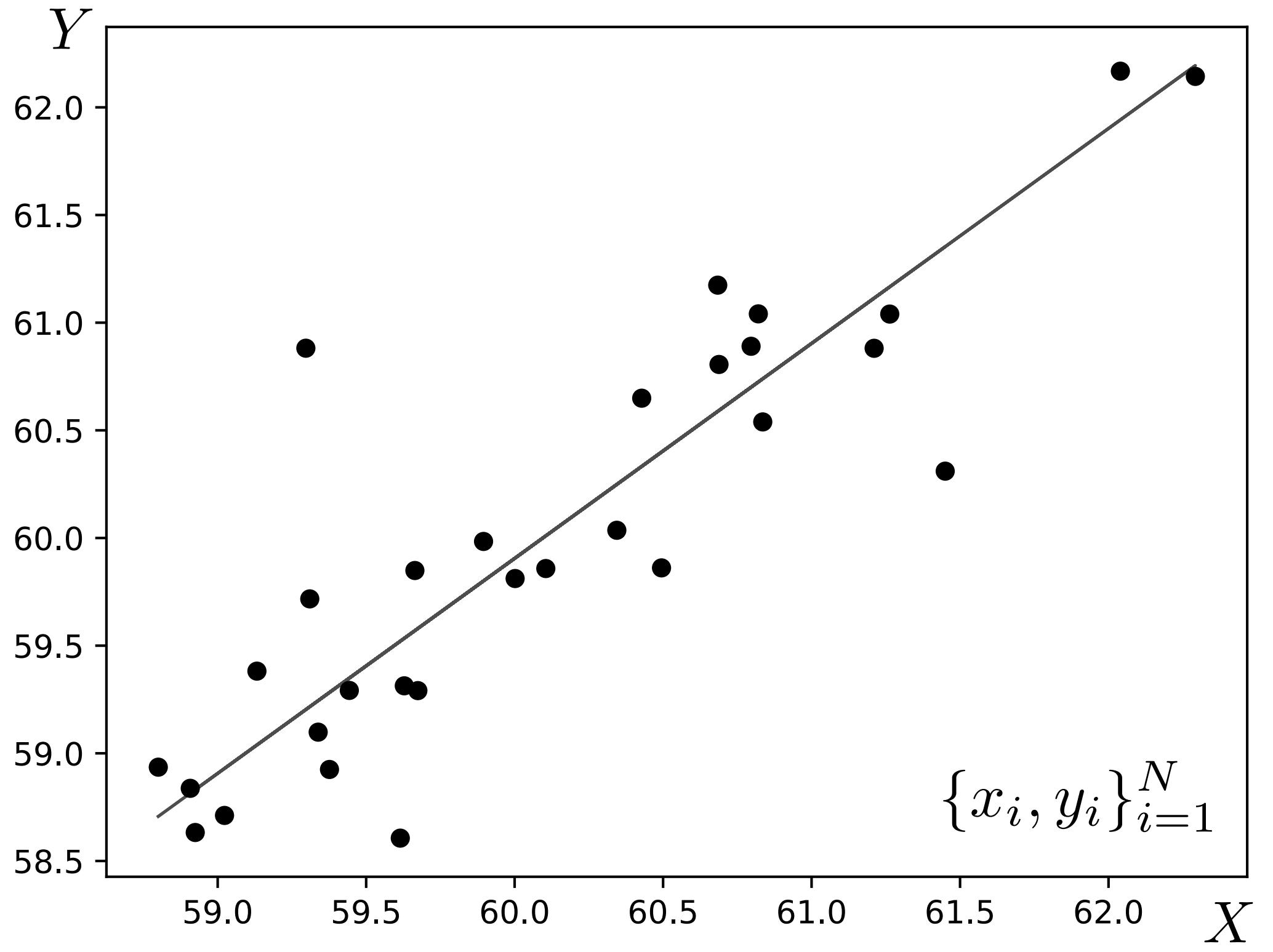
Training data



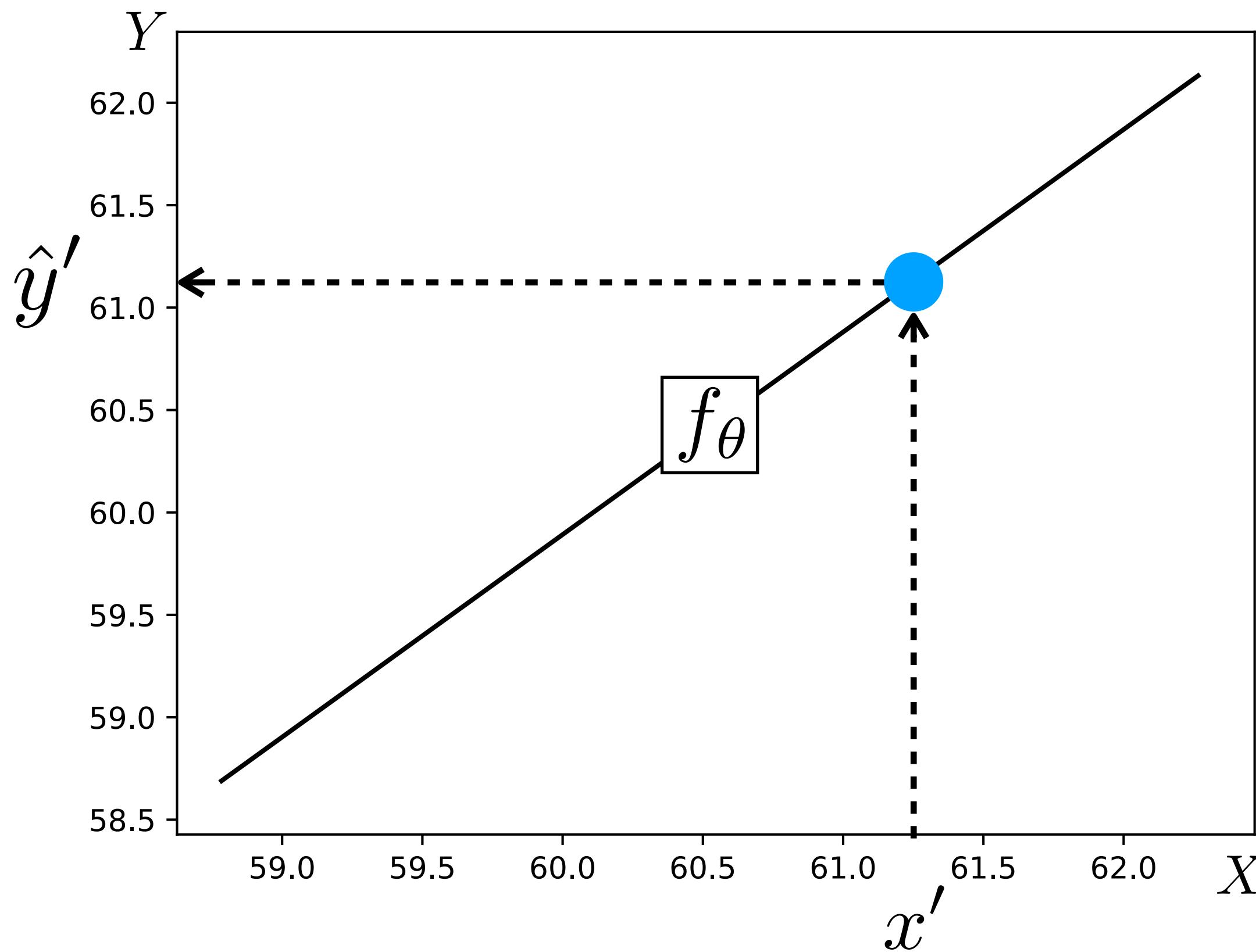
Test query



Training data



Test query



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

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:

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

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

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

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

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

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

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

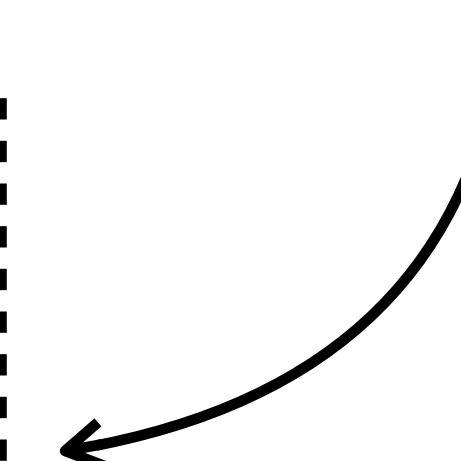
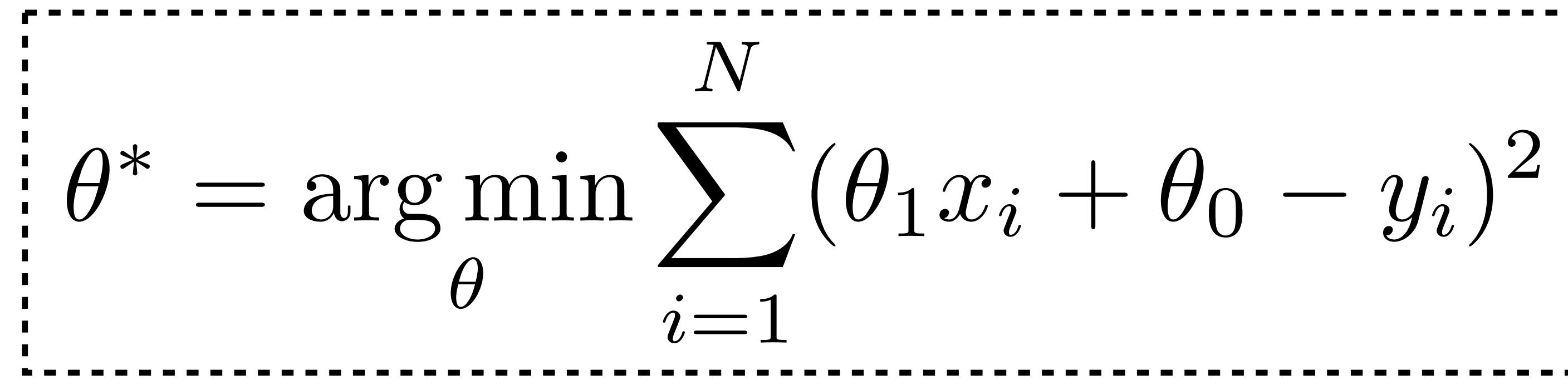
$$\boxed{\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}$$

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

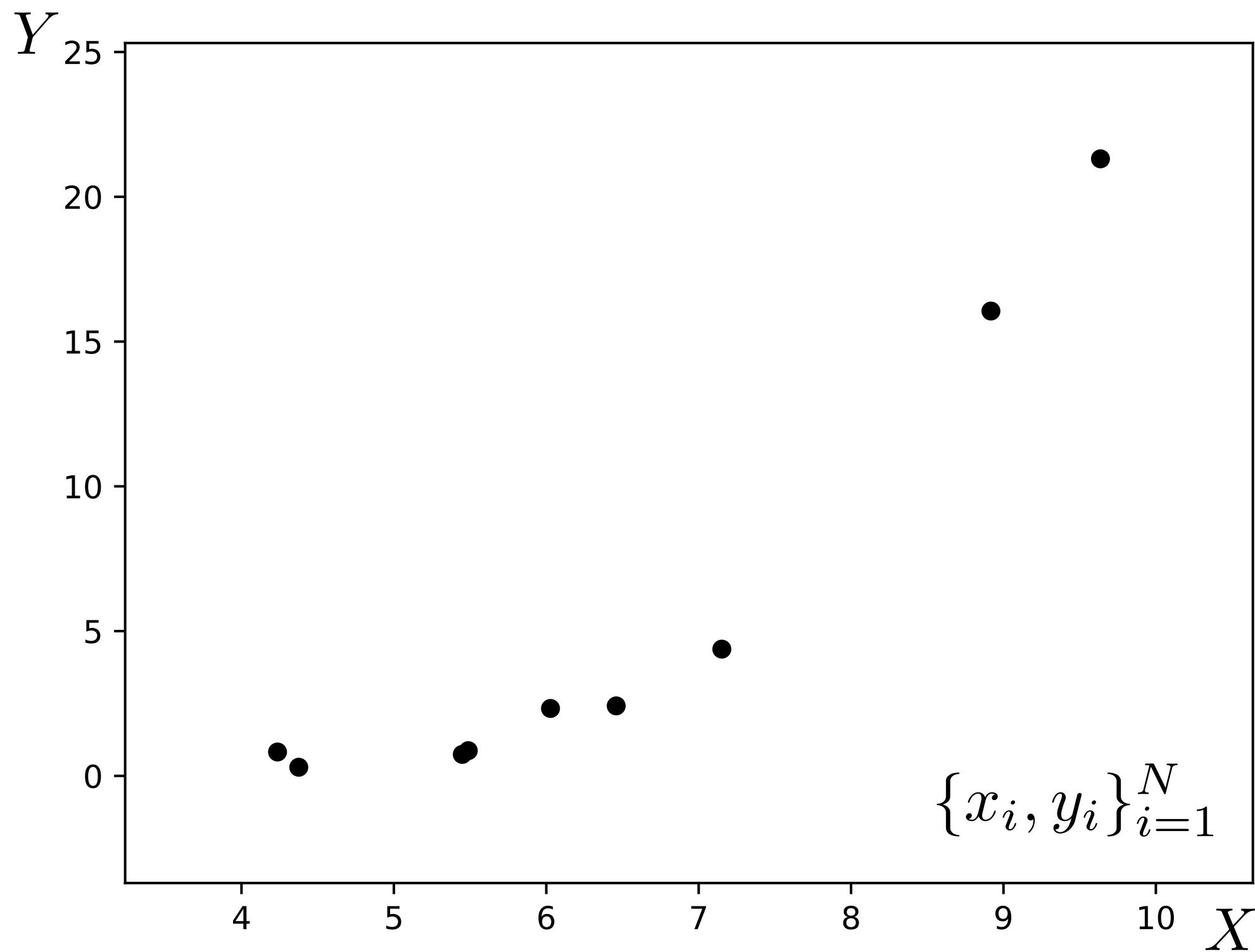
Training data

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 arrows point to different parts of the formula: one from 'Hypothesis space' to the f in the argument of the minimum, one from 'Training data' to the summation symbol, and one from 'Objective function (loss)' to the loss function \mathcal{L} .

The Problem of Generalization

Linear regression

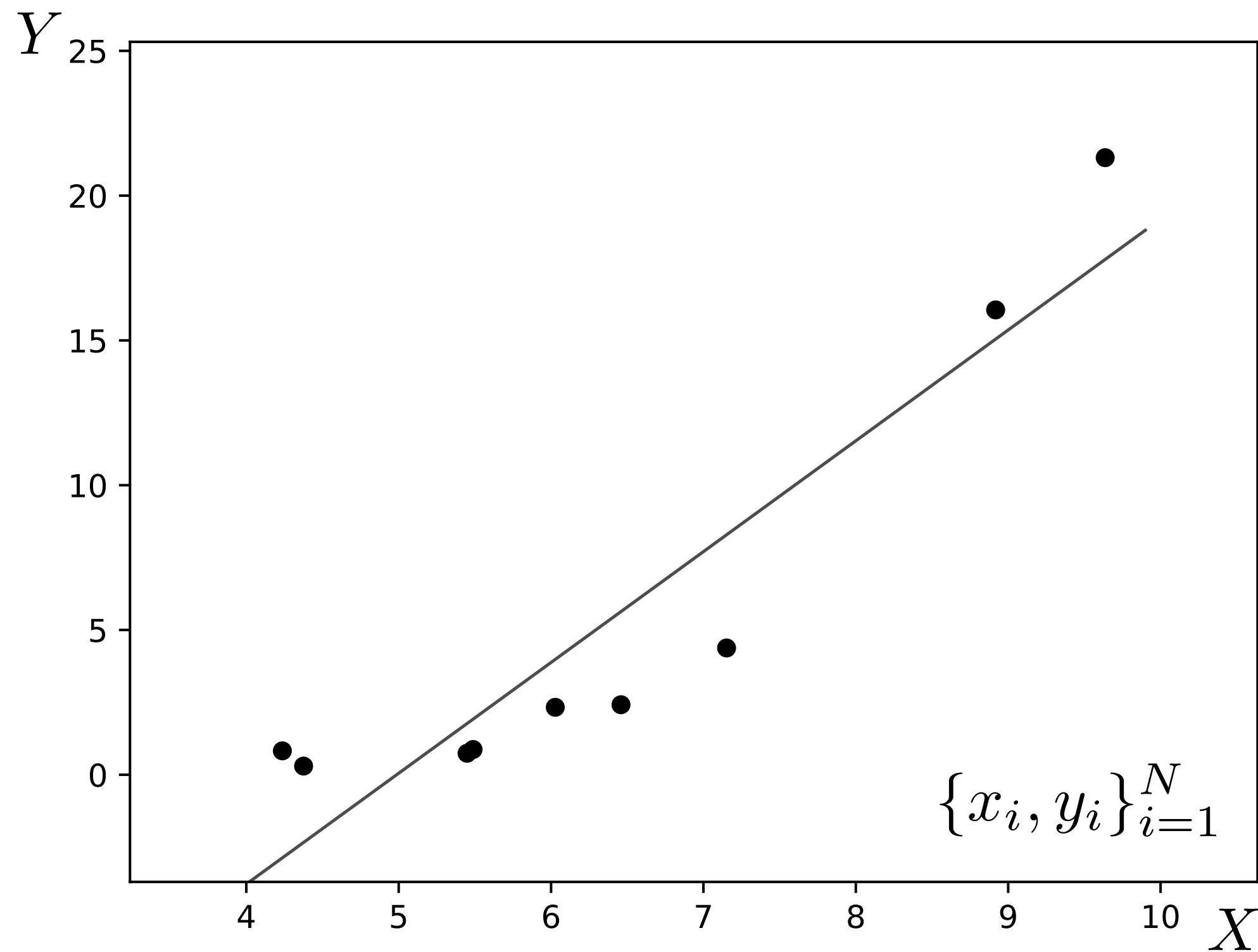
Training data



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

Linear regression

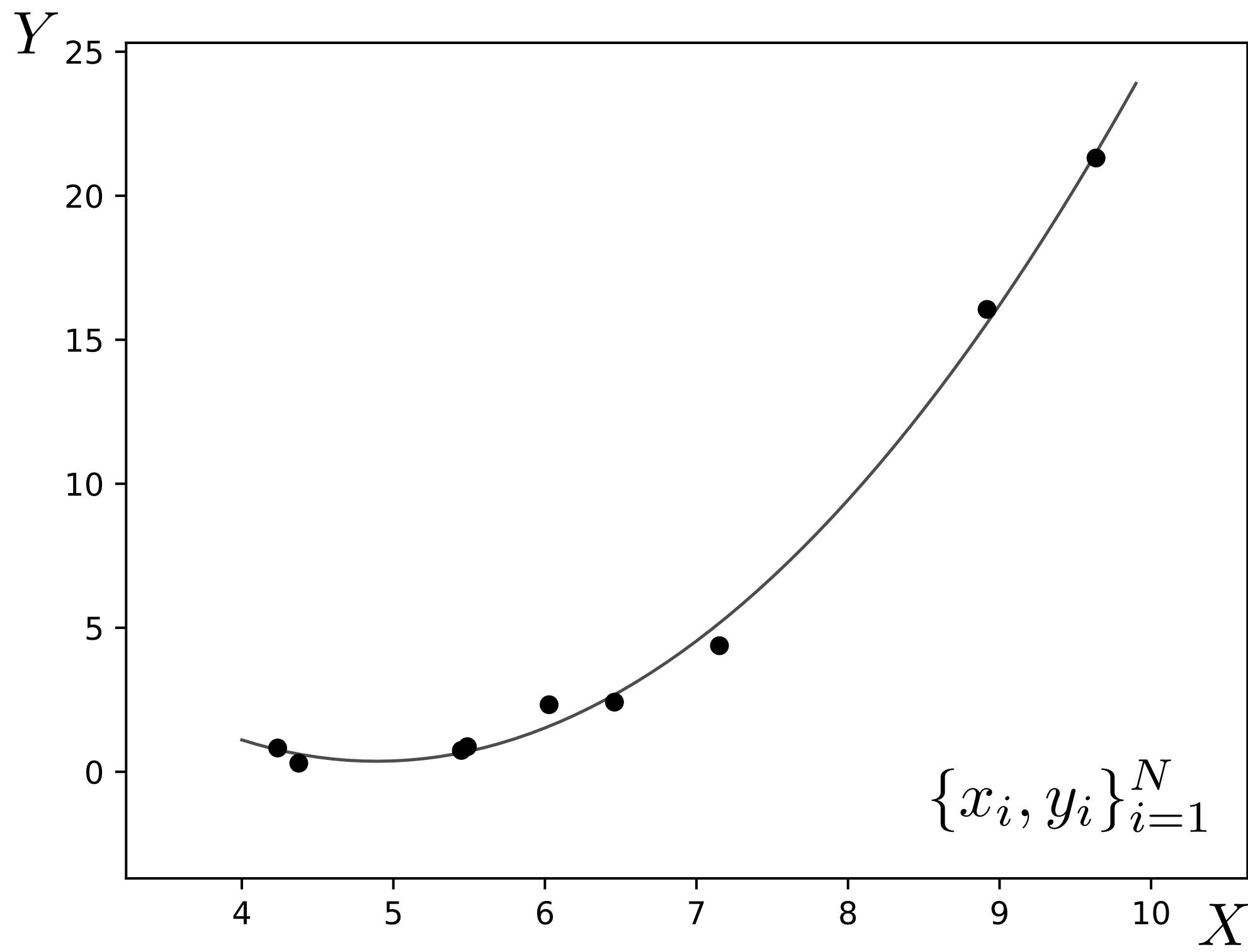
Training data



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

Polynomial regression

Training data

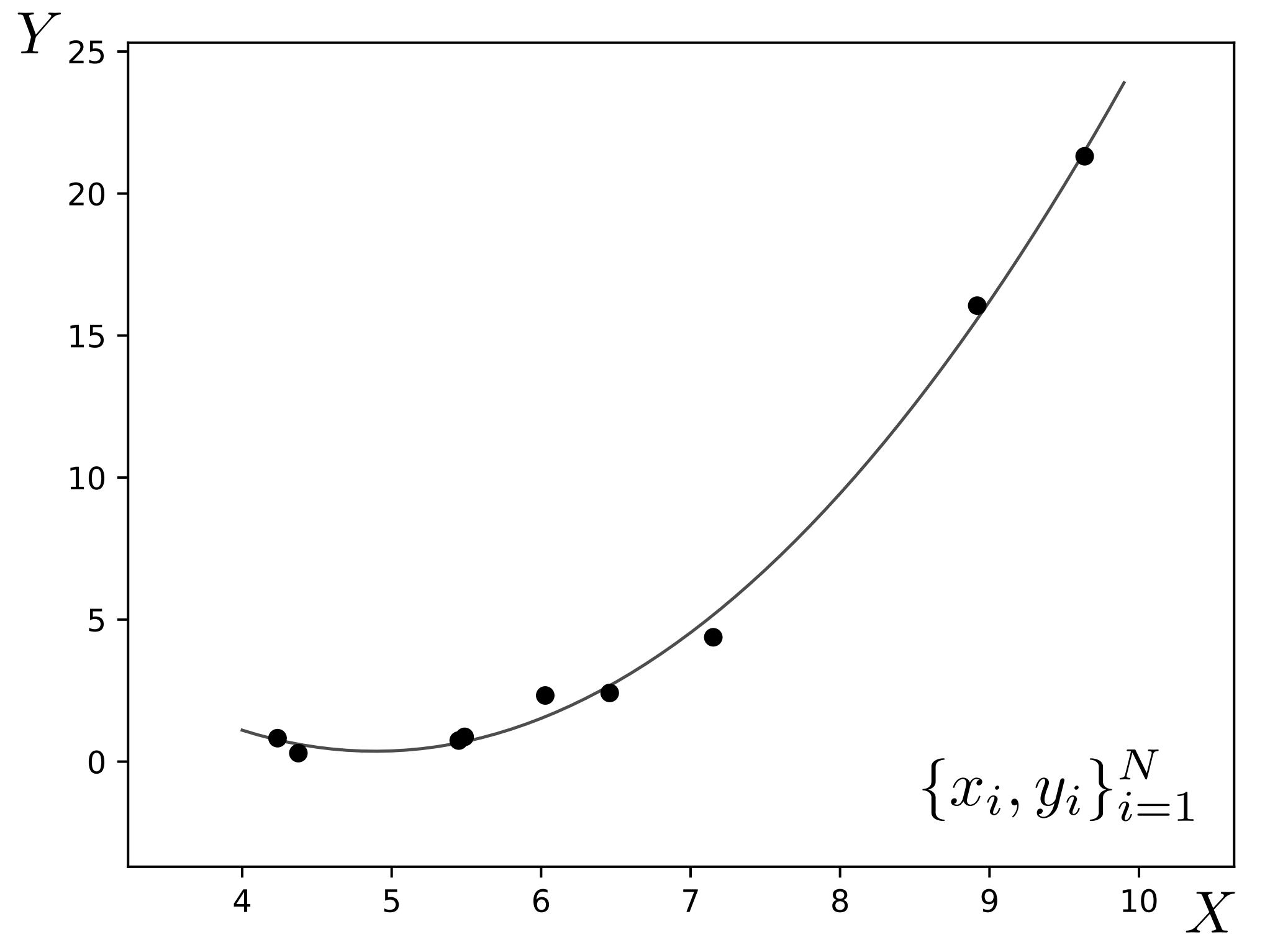


$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

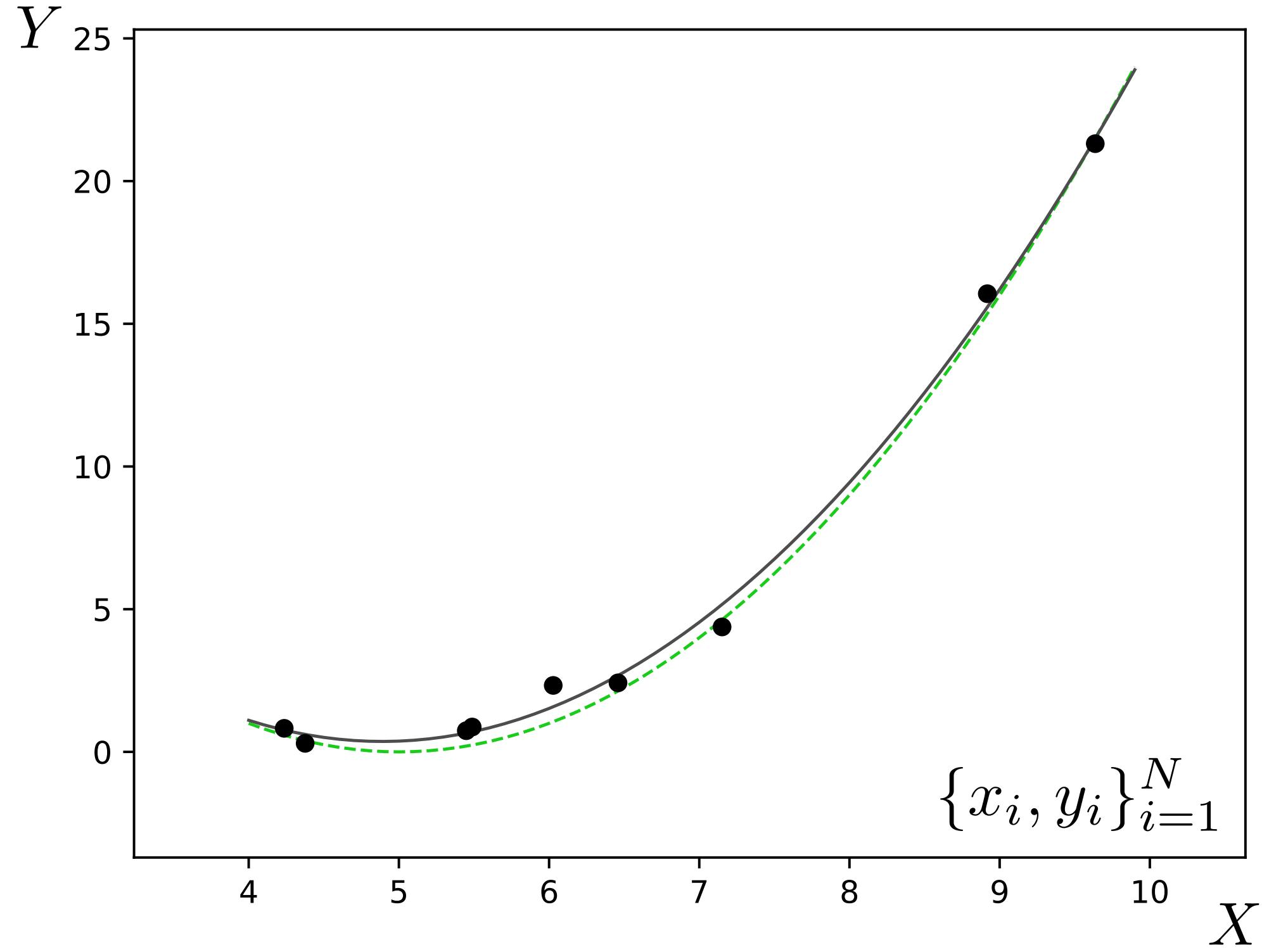
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

K-th degree polynomial regression

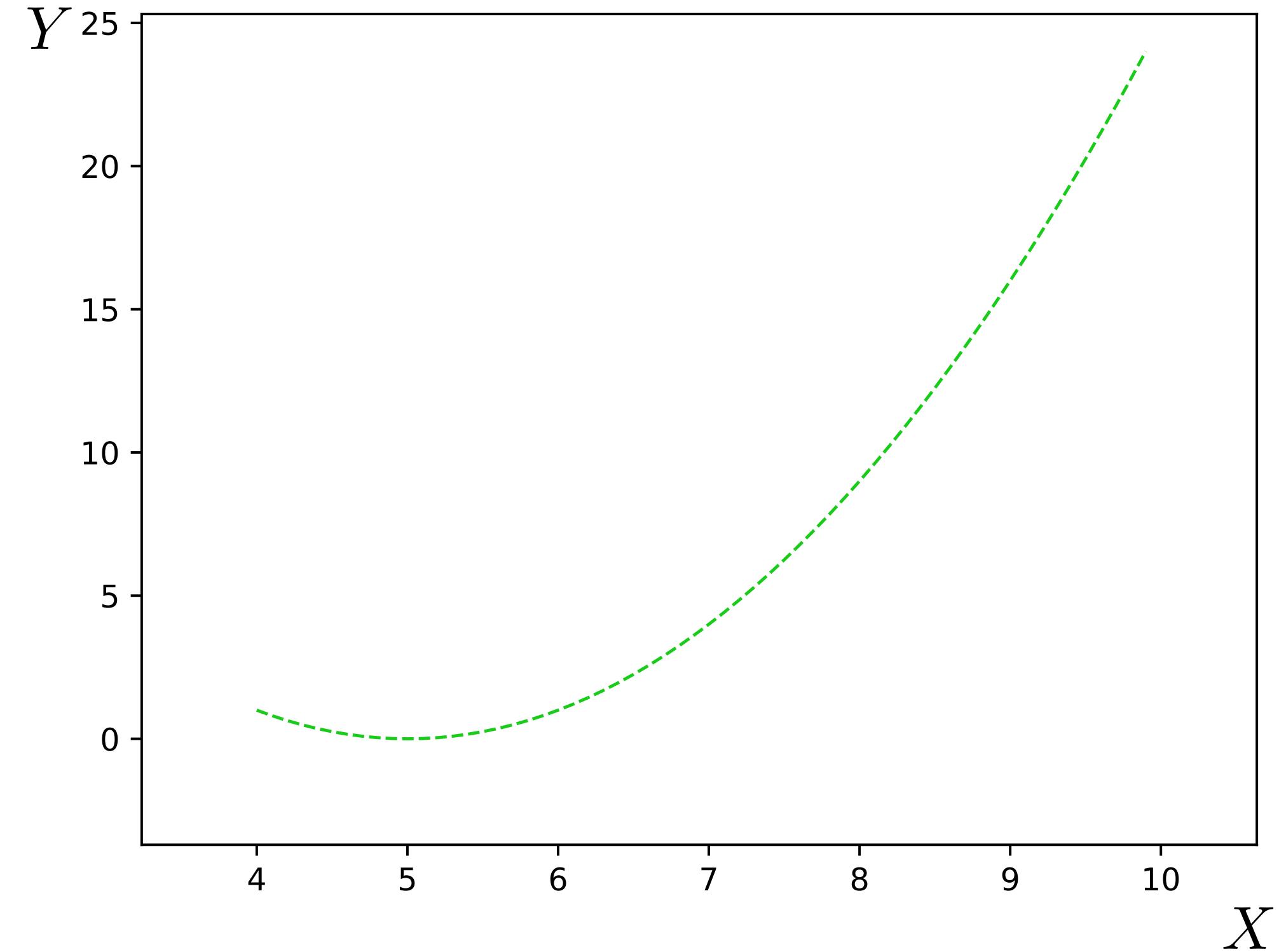
Training data



Training data



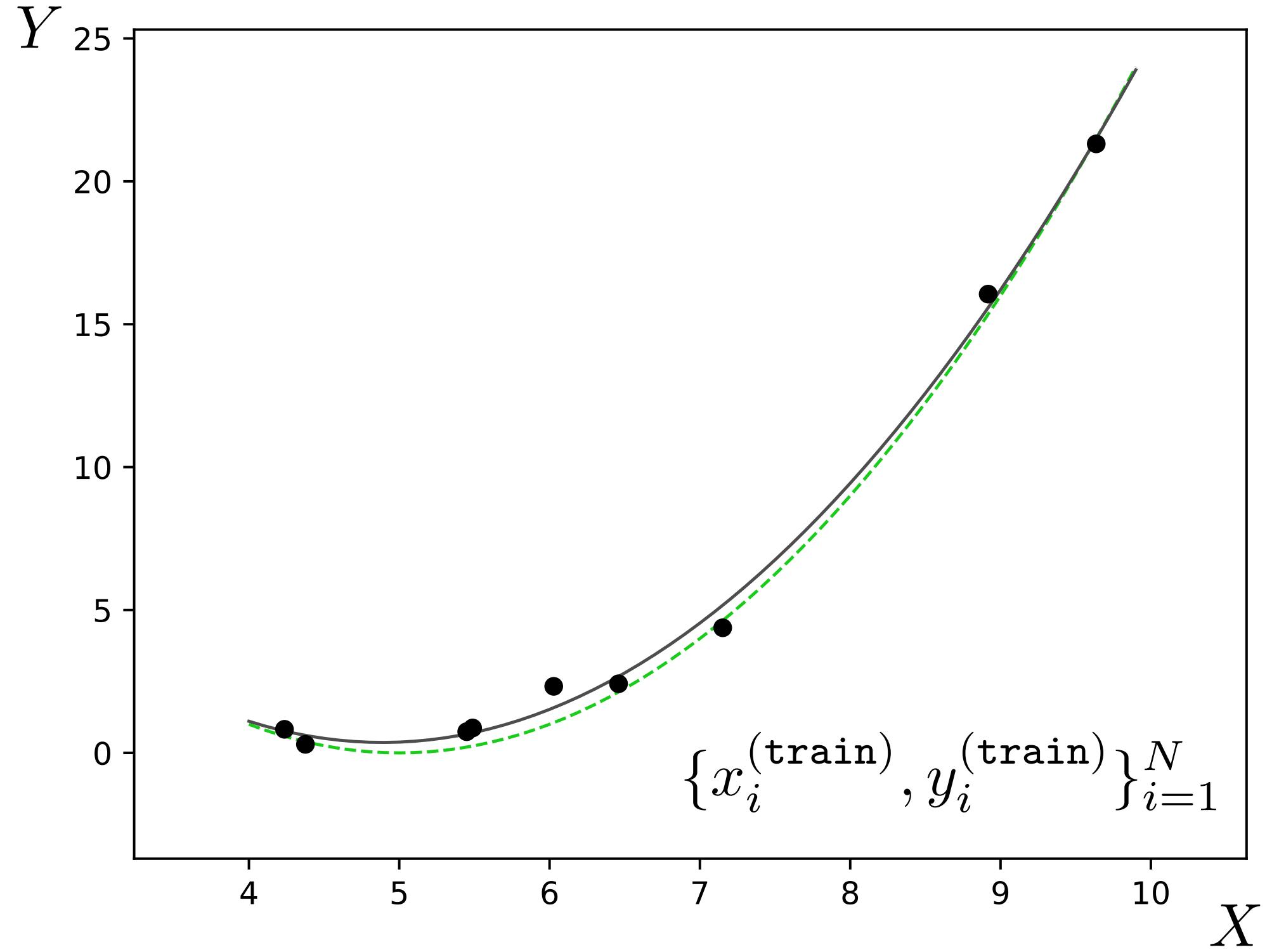
Test data



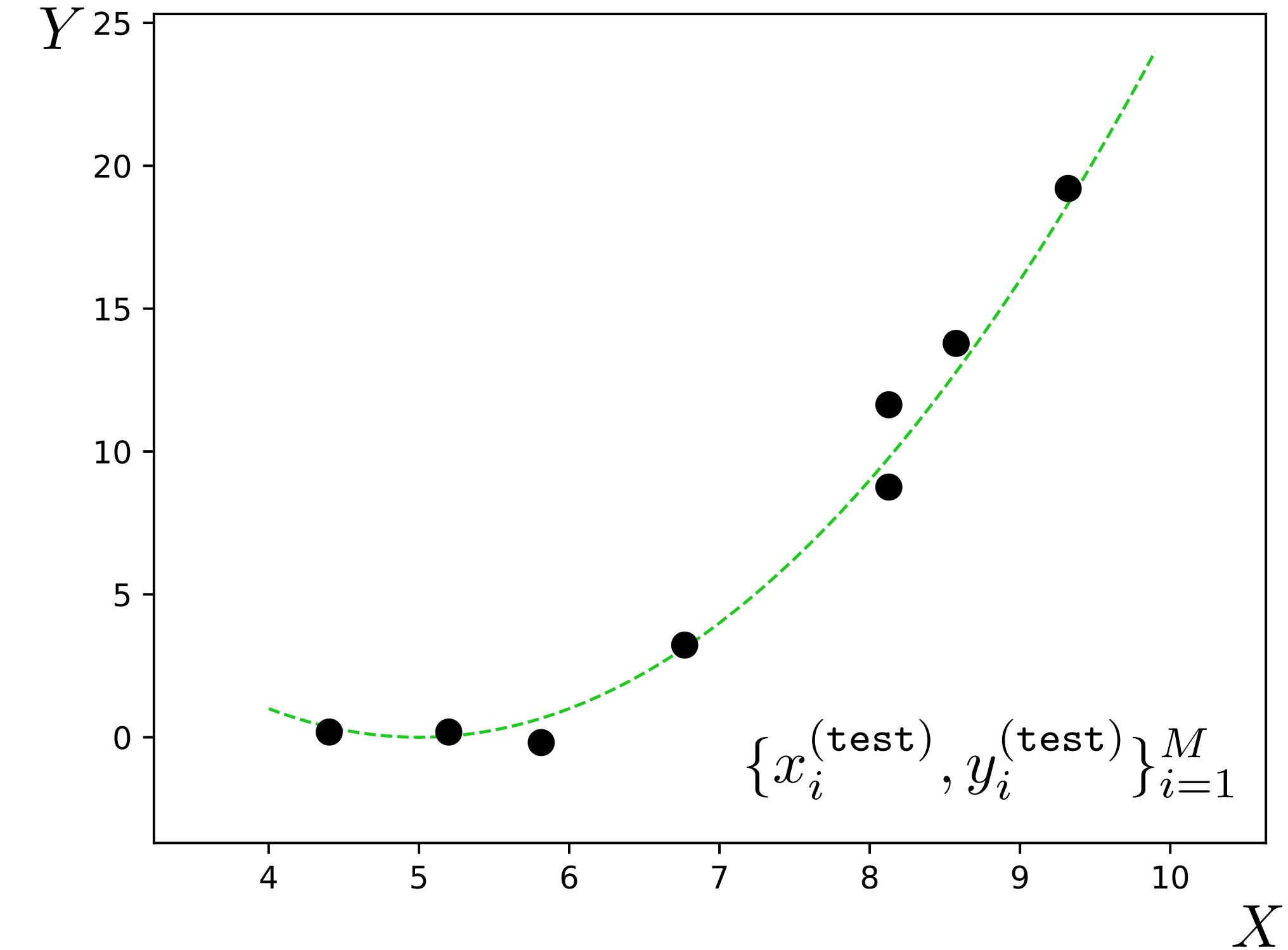
True data-generating process

p_{data}

Training data



Test data



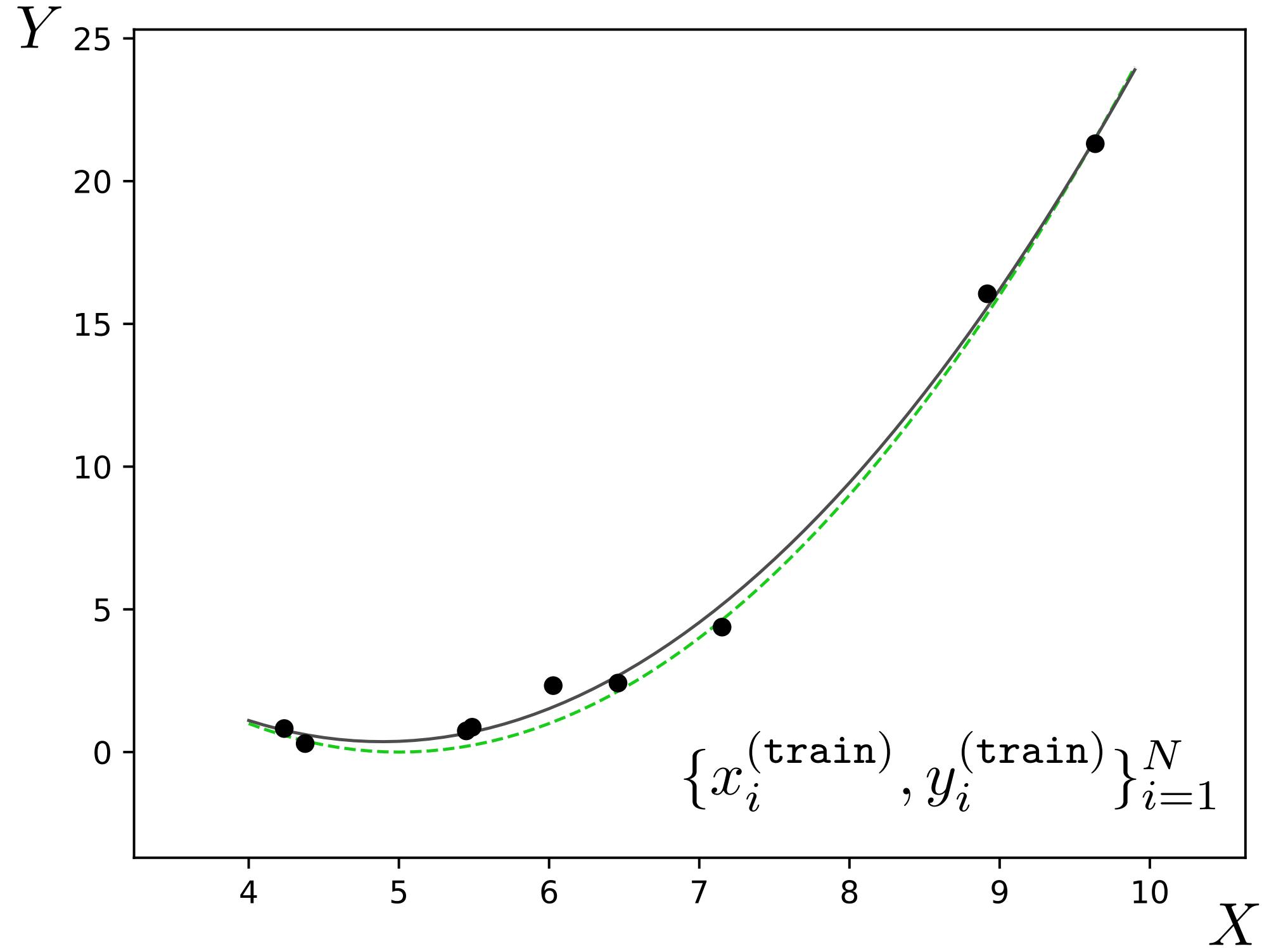
True data-generating process

p_{data}

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

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

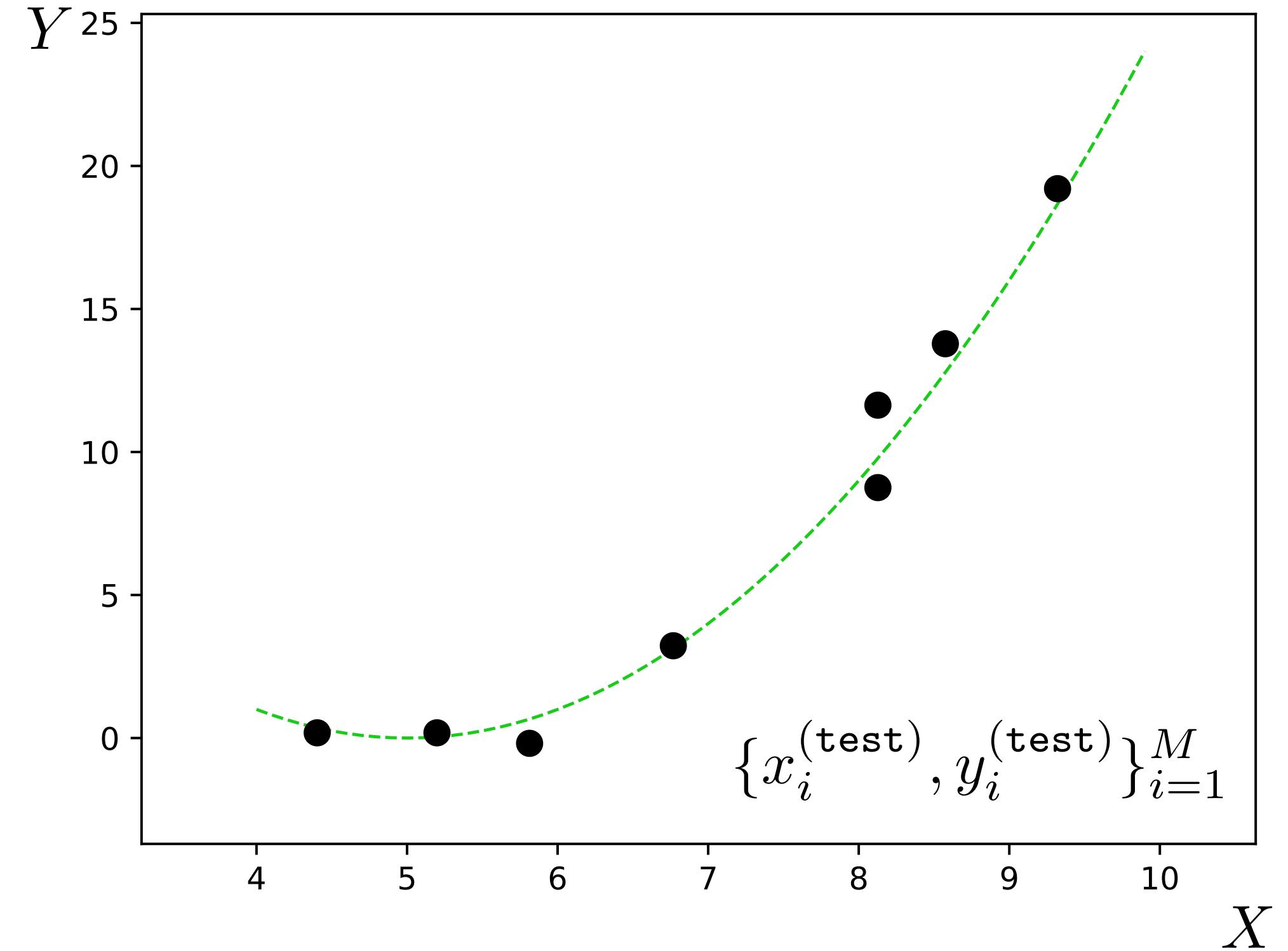
Training data



Training objective:

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

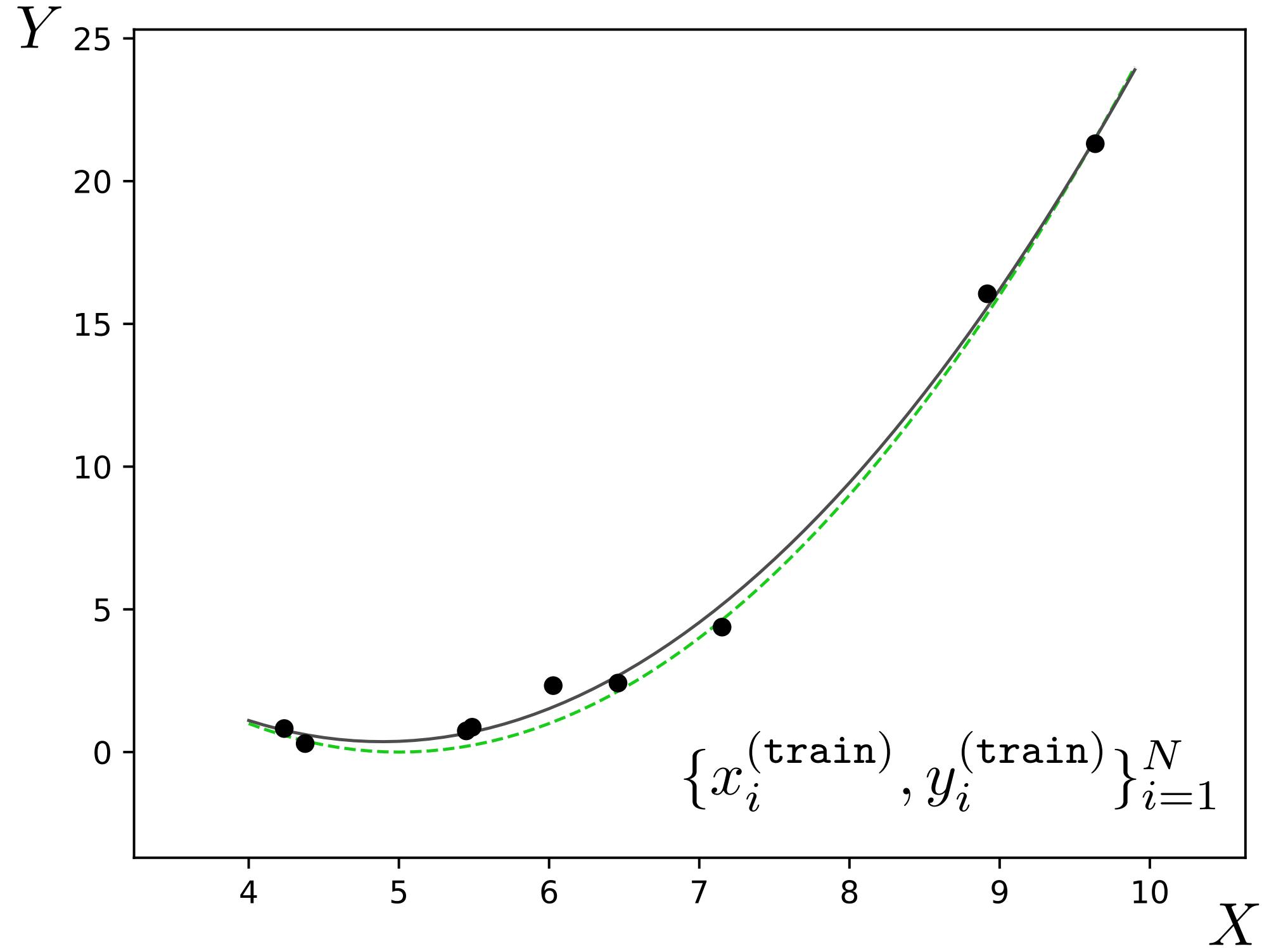
Test data



Test time evaluation:

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

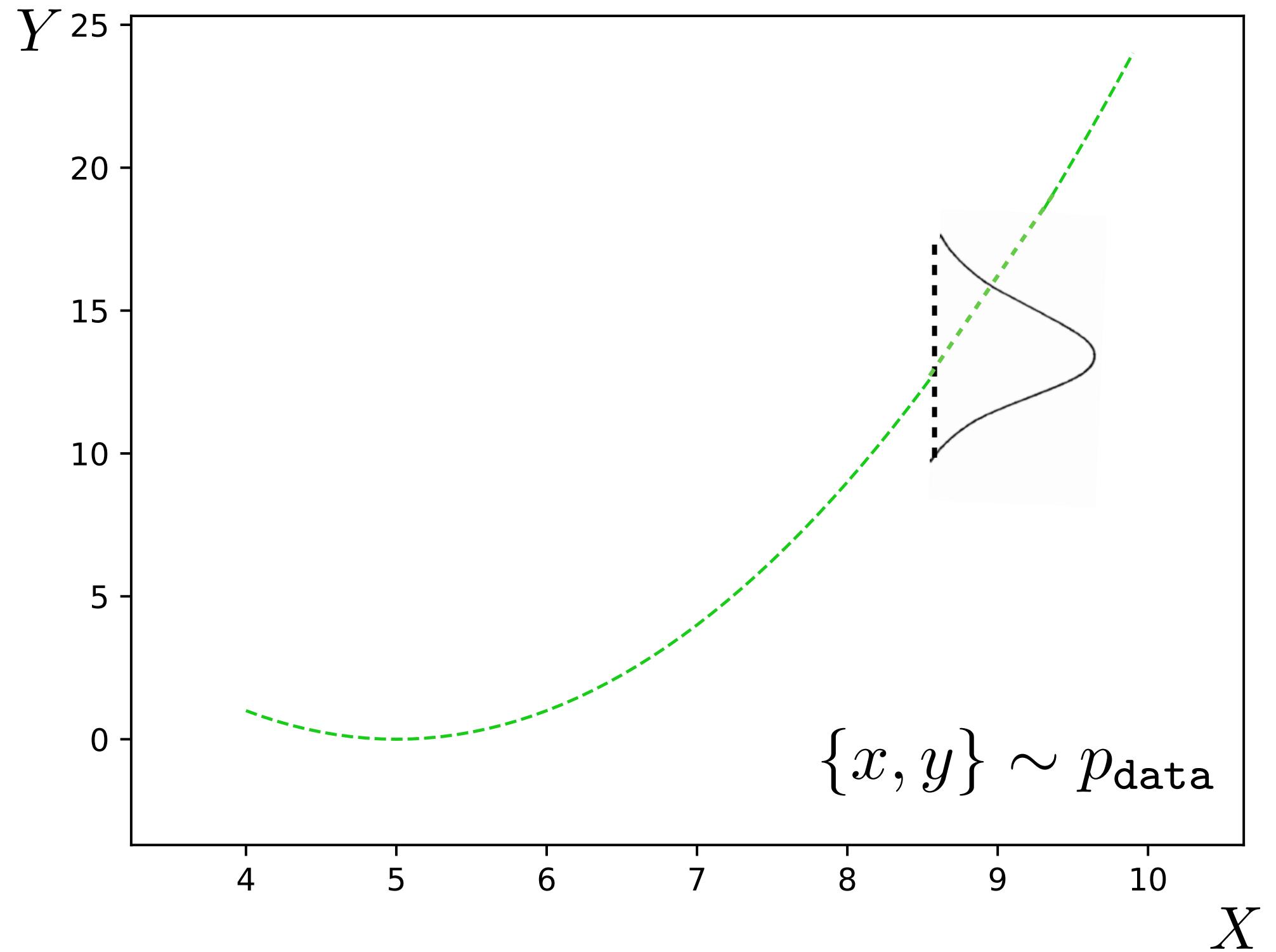
Training data



Training objective:

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

Test data

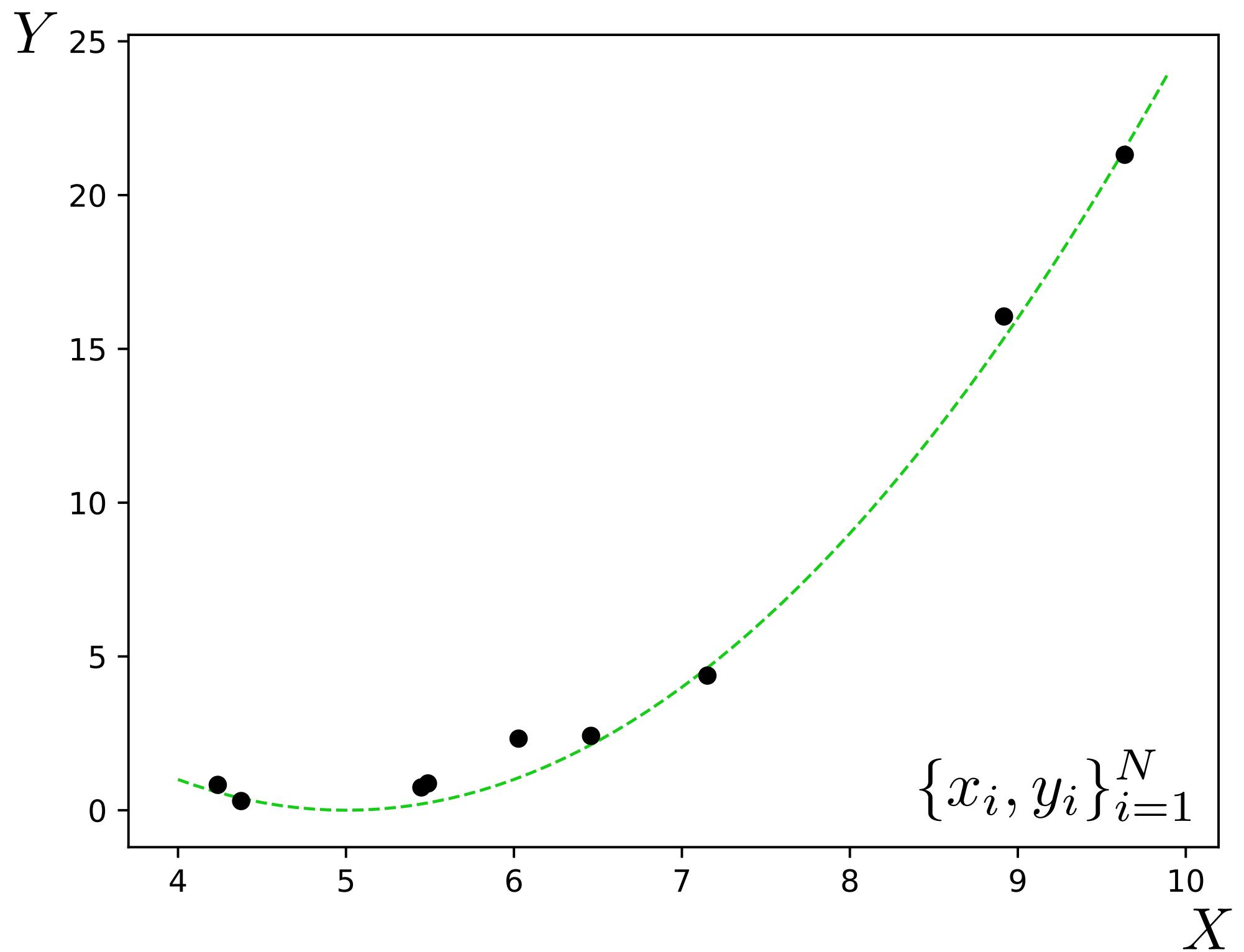


True objective:

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

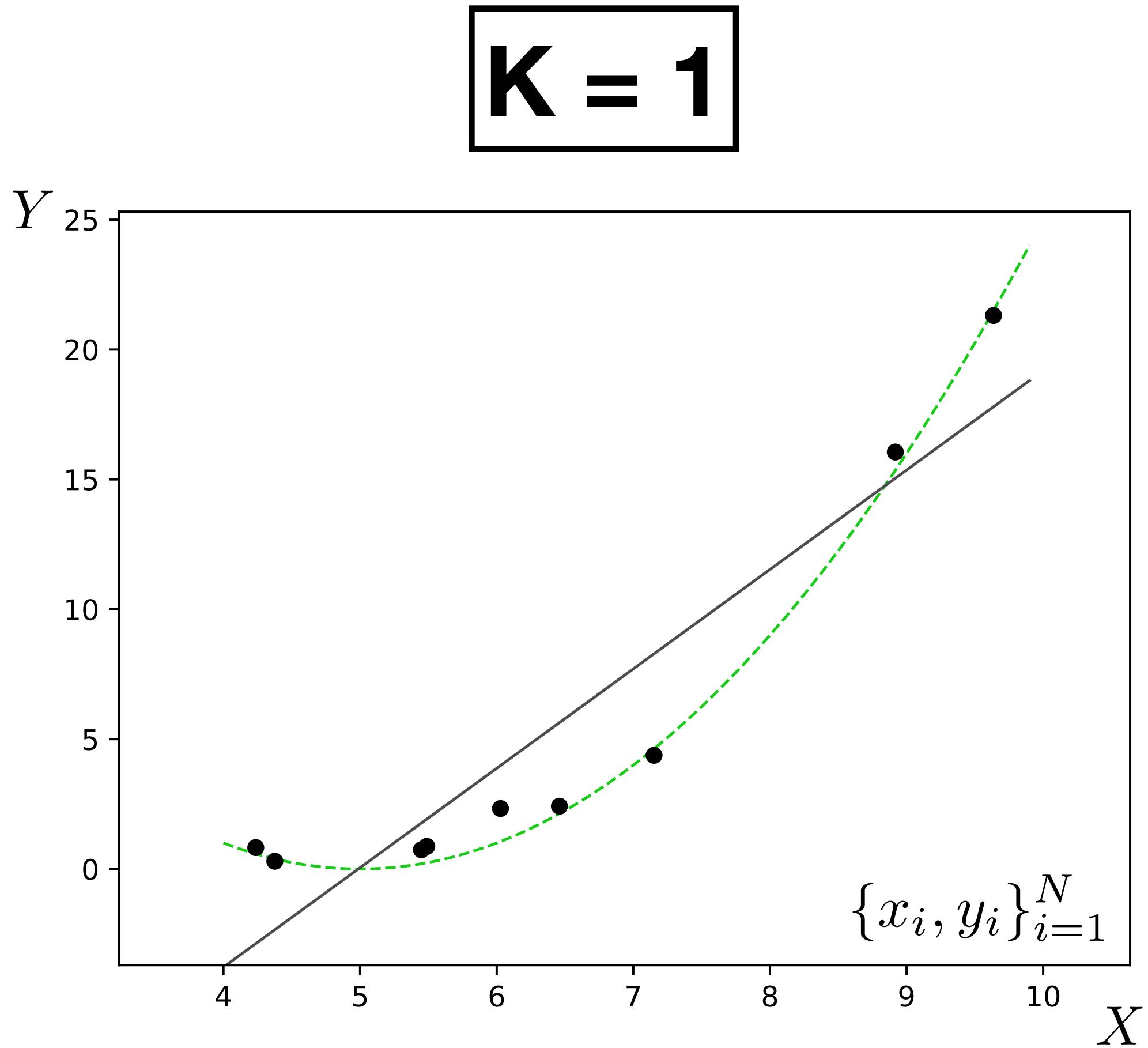
What happens as we add more basis functions?

Training data



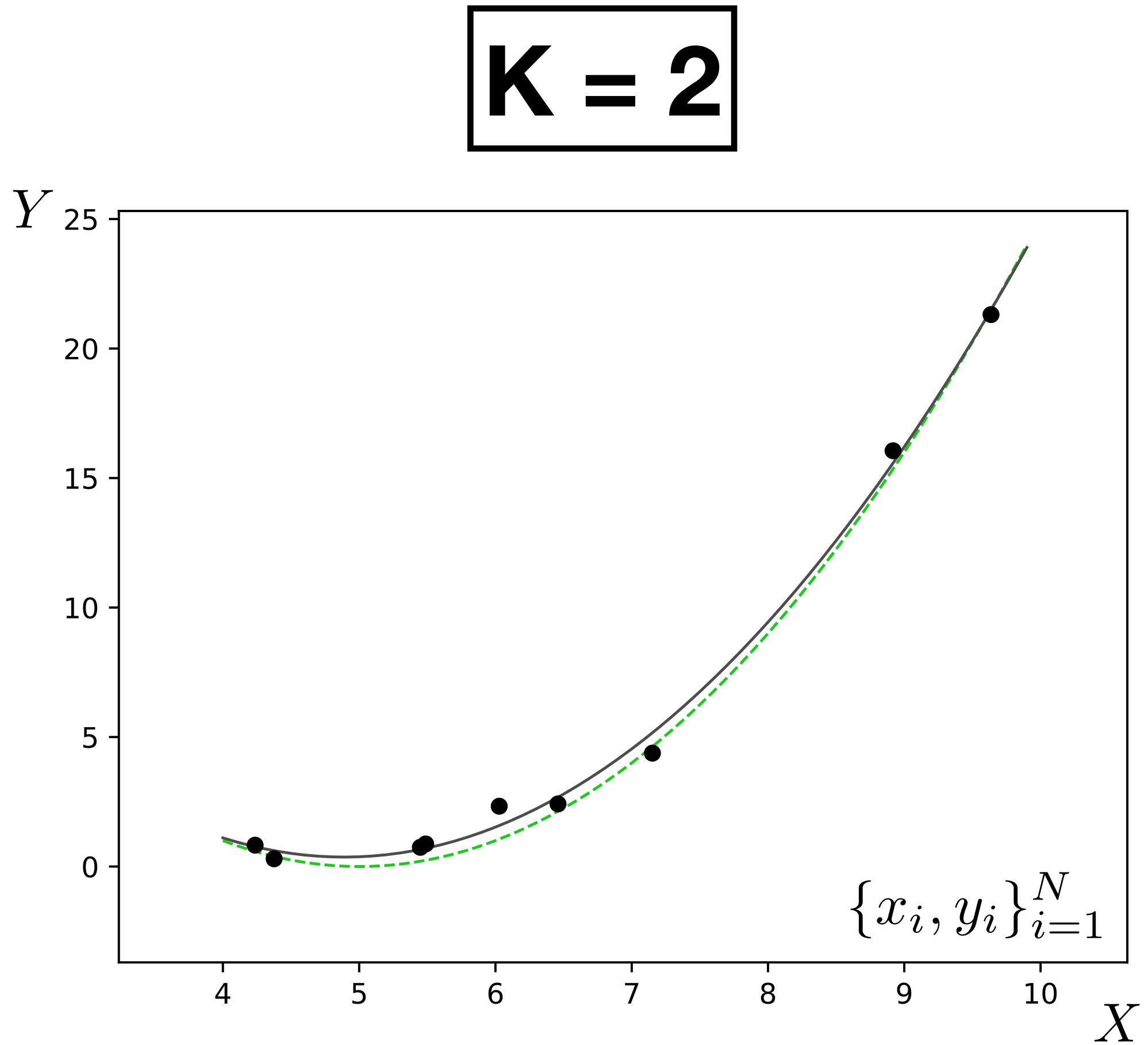
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?



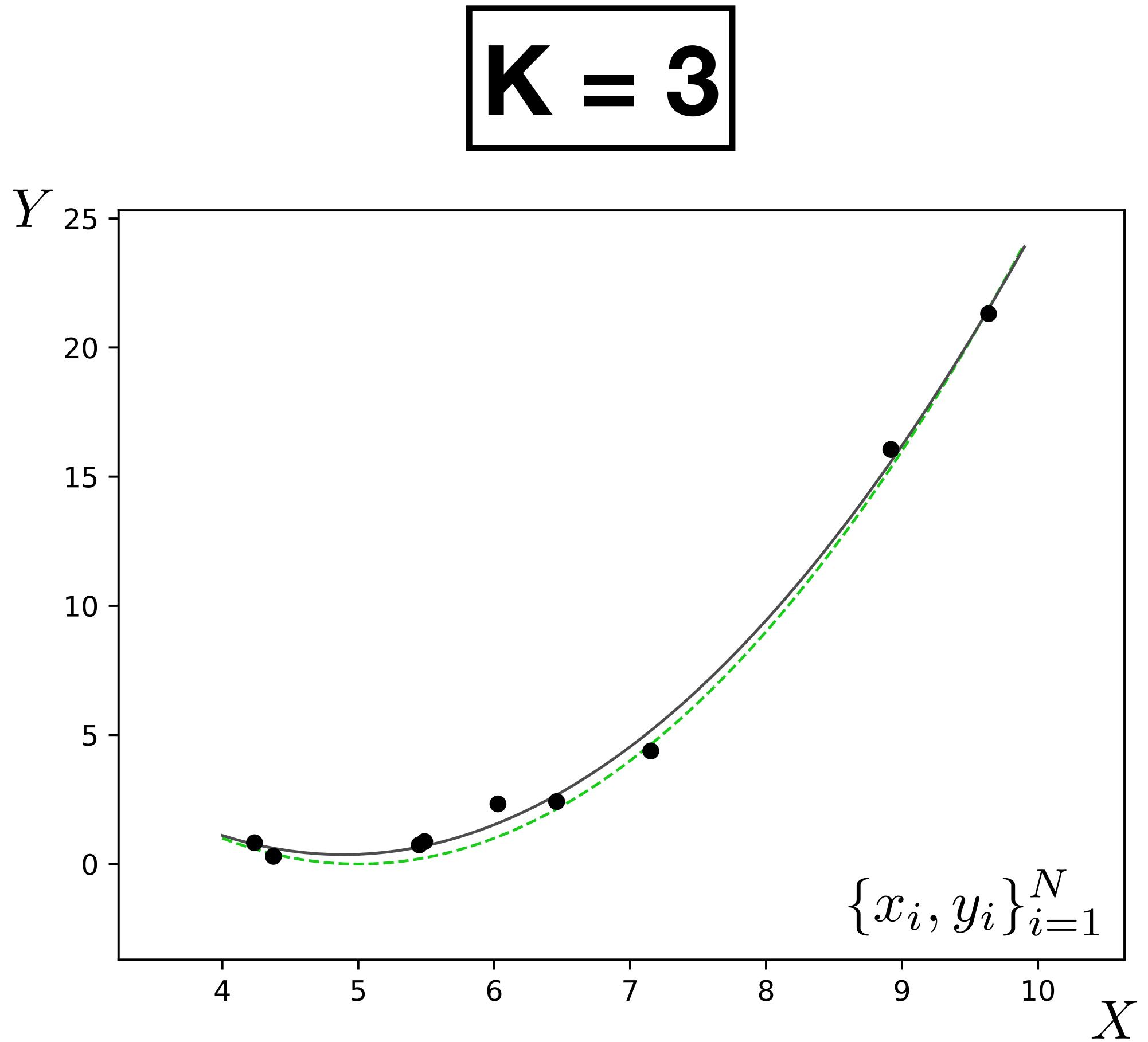
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



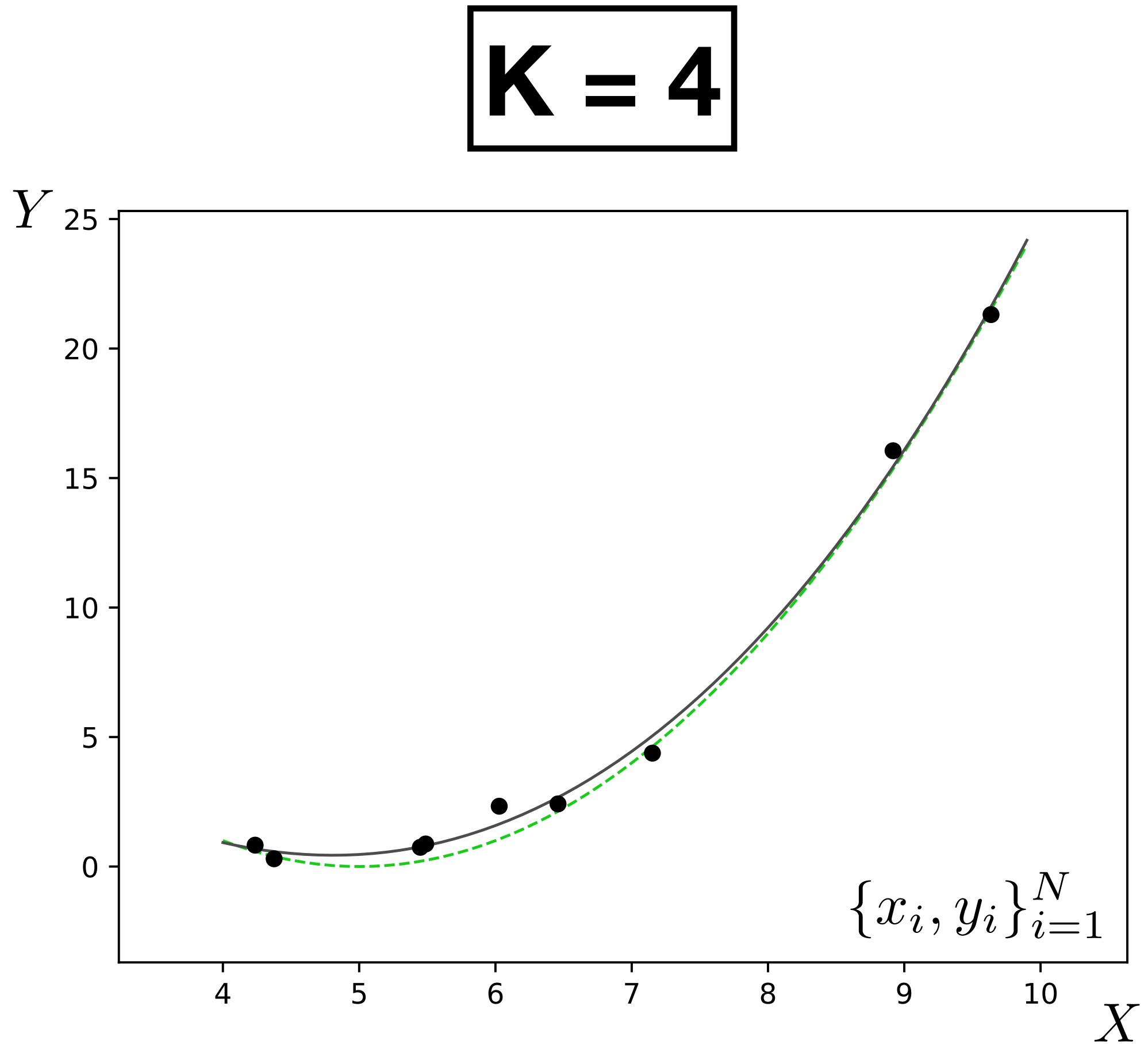
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



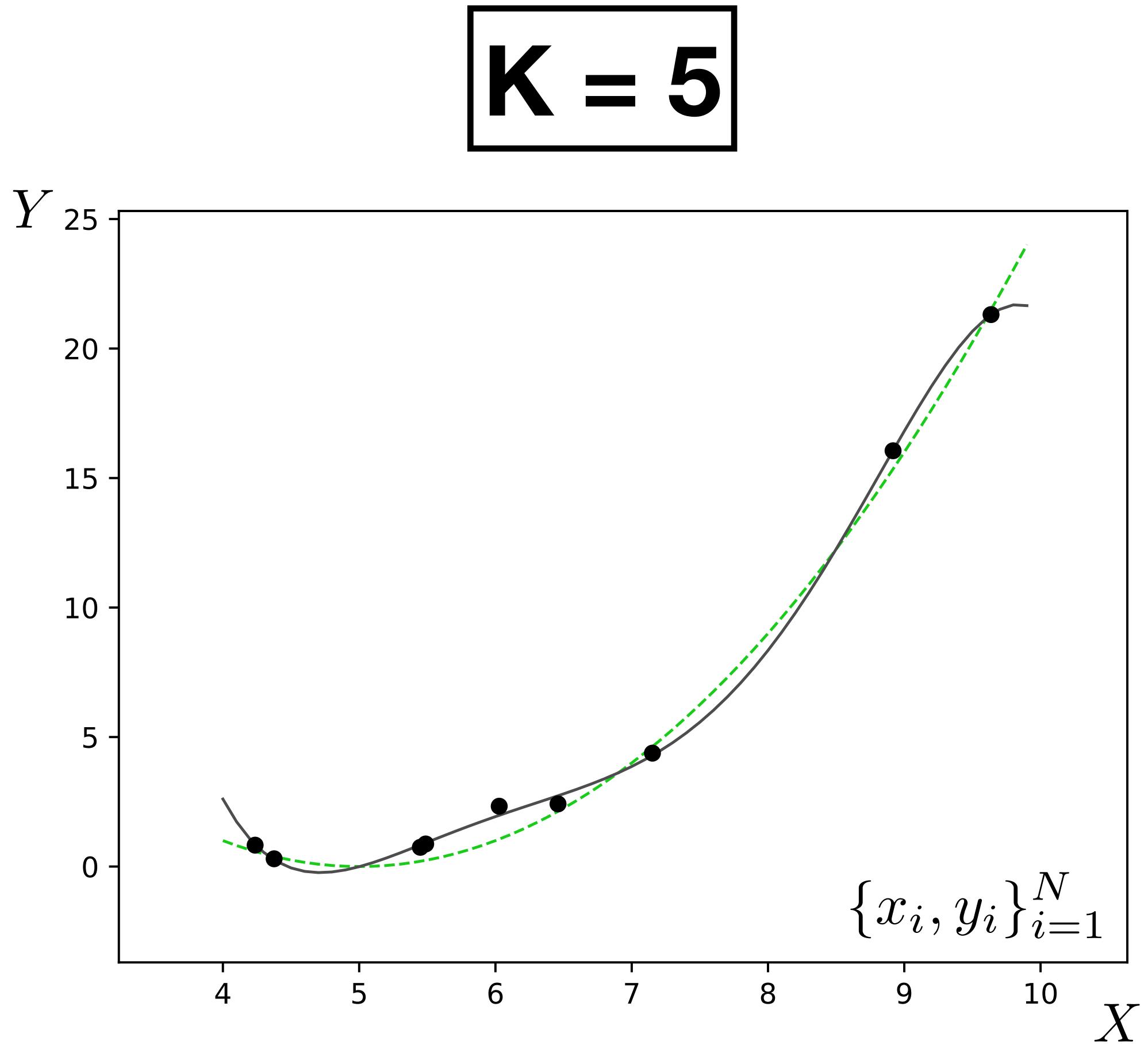
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



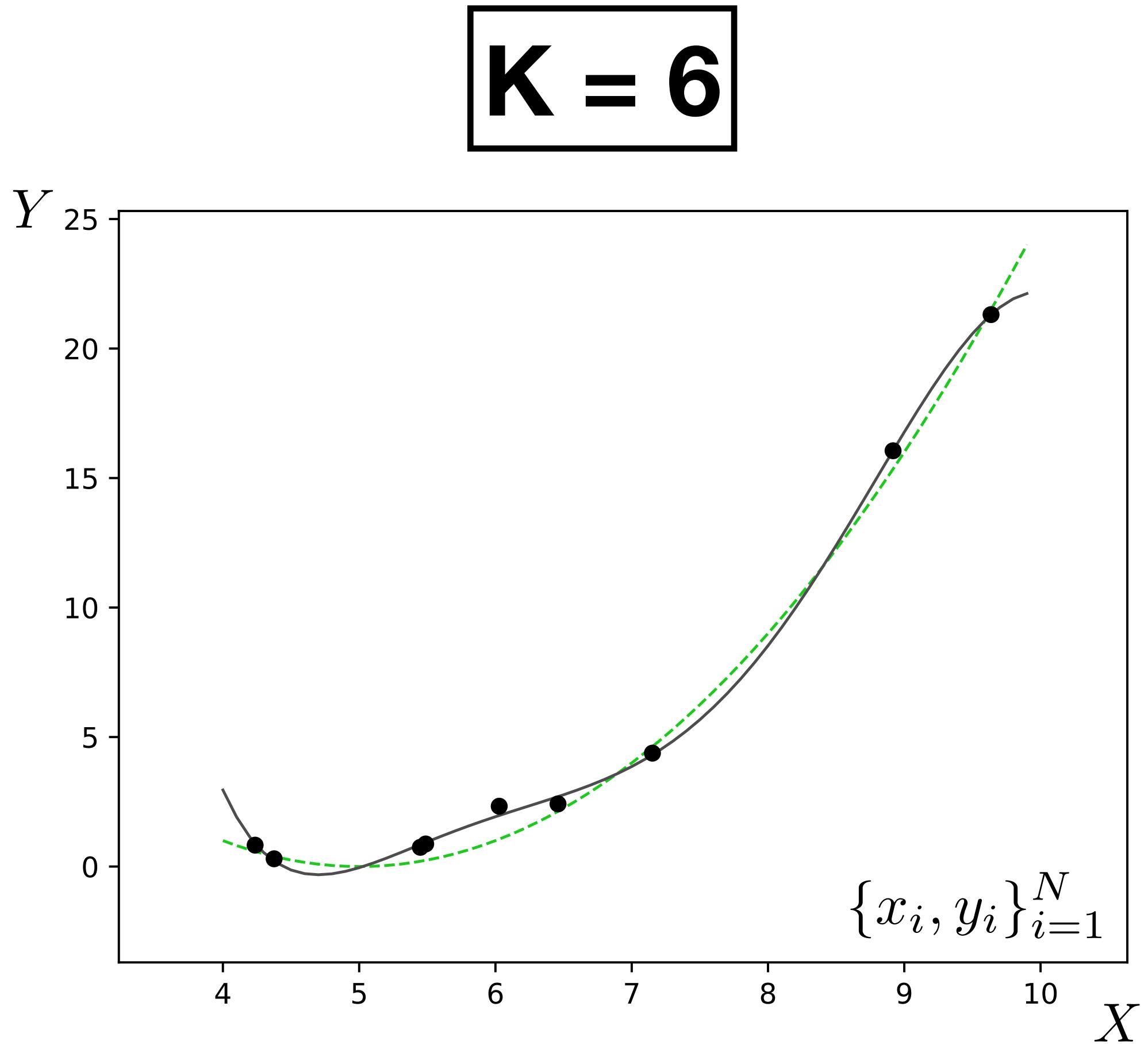
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



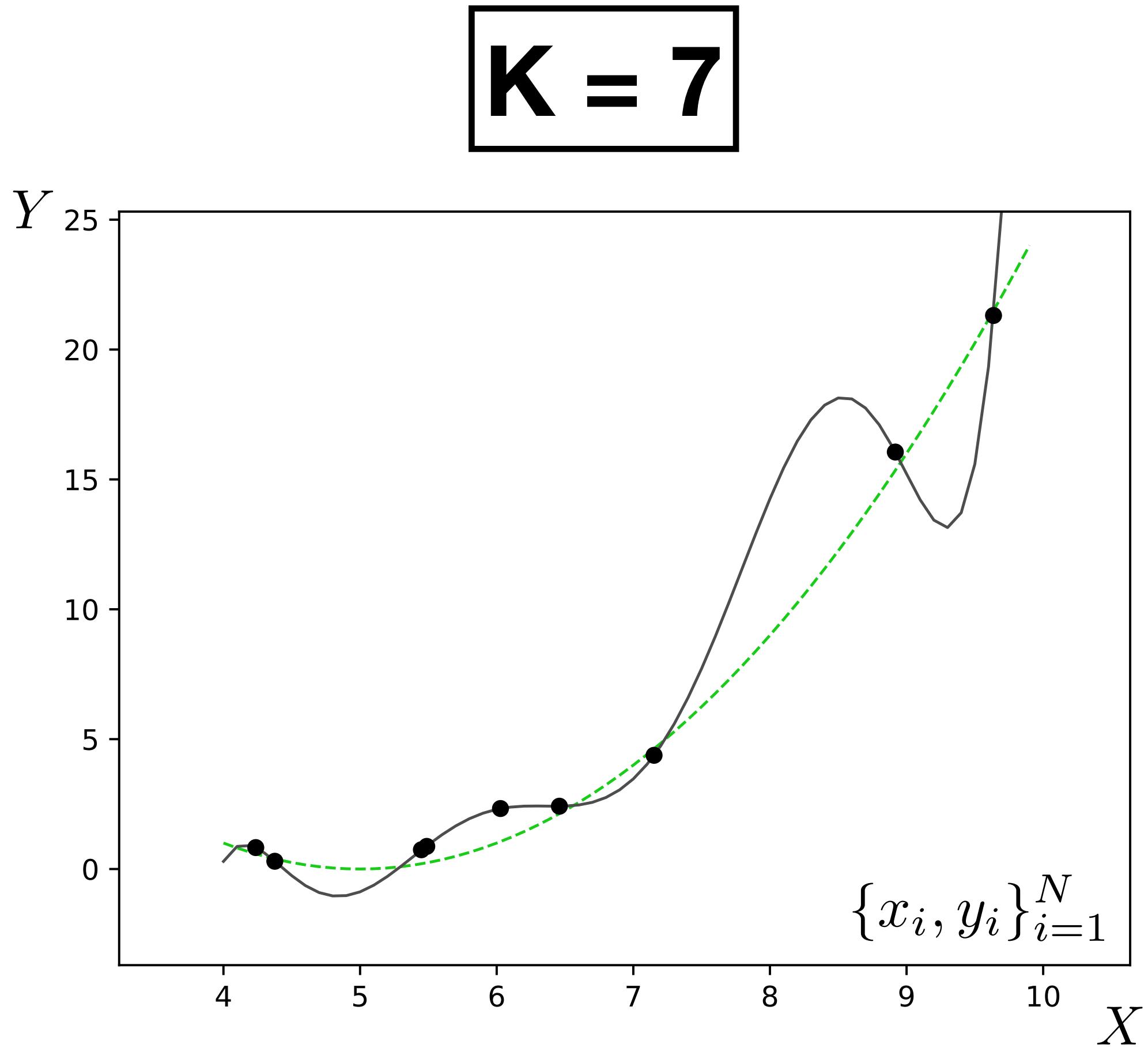
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



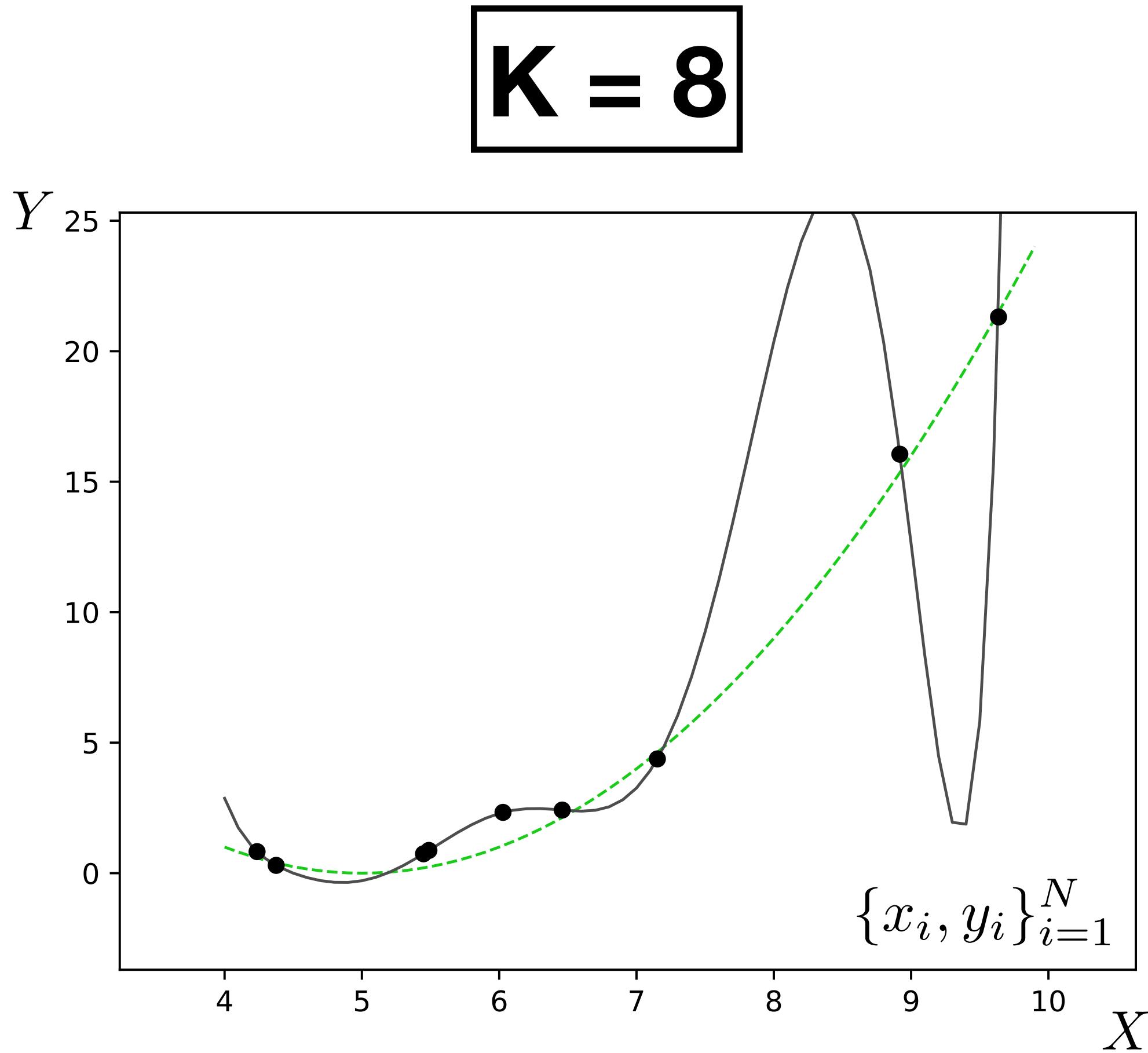
$$f_\theta(x) = \sum_{k=0}^K \theta_k x^k$$

What happens as we add more basis functions?



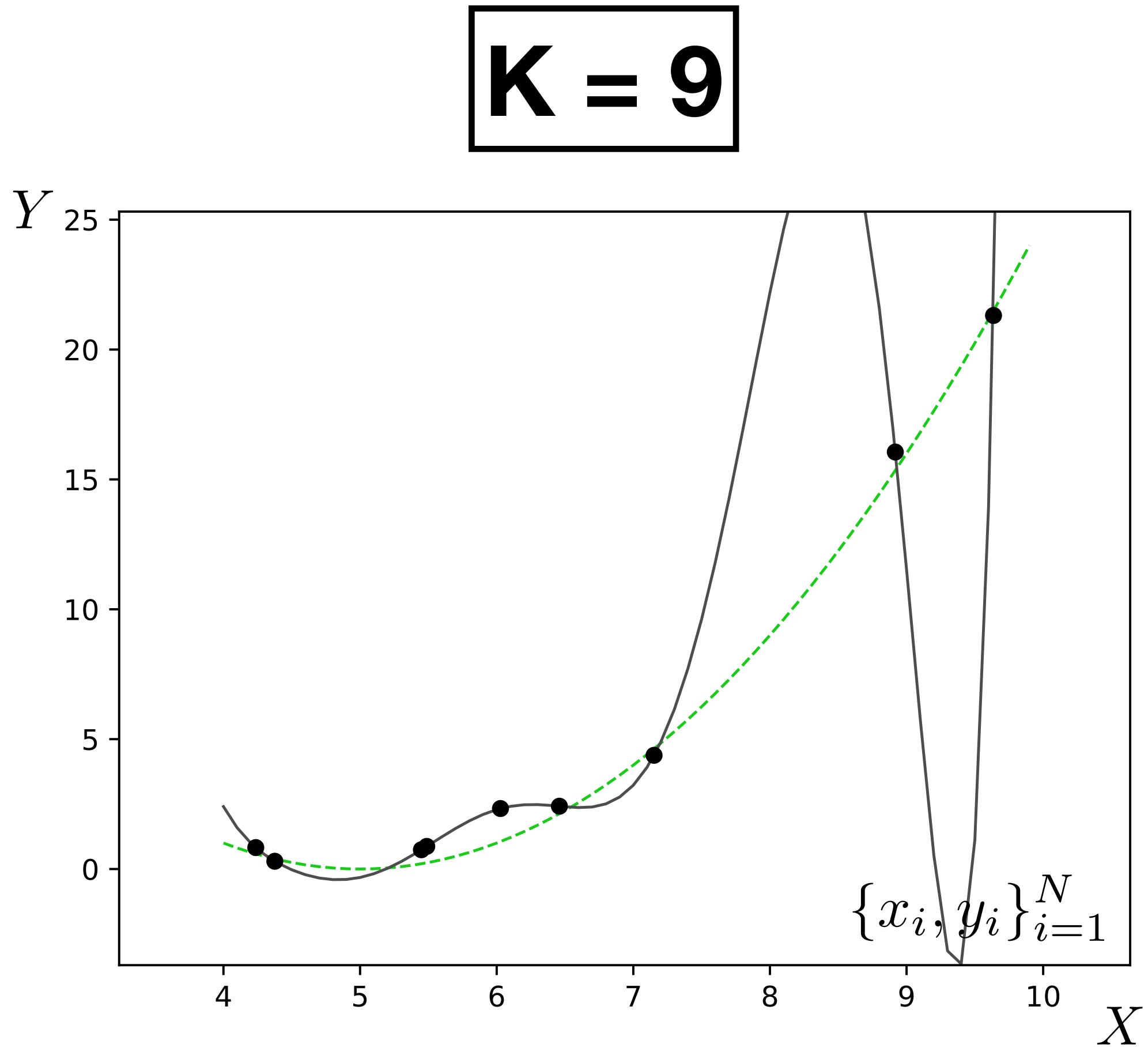
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

What happens as we add more basis functions?



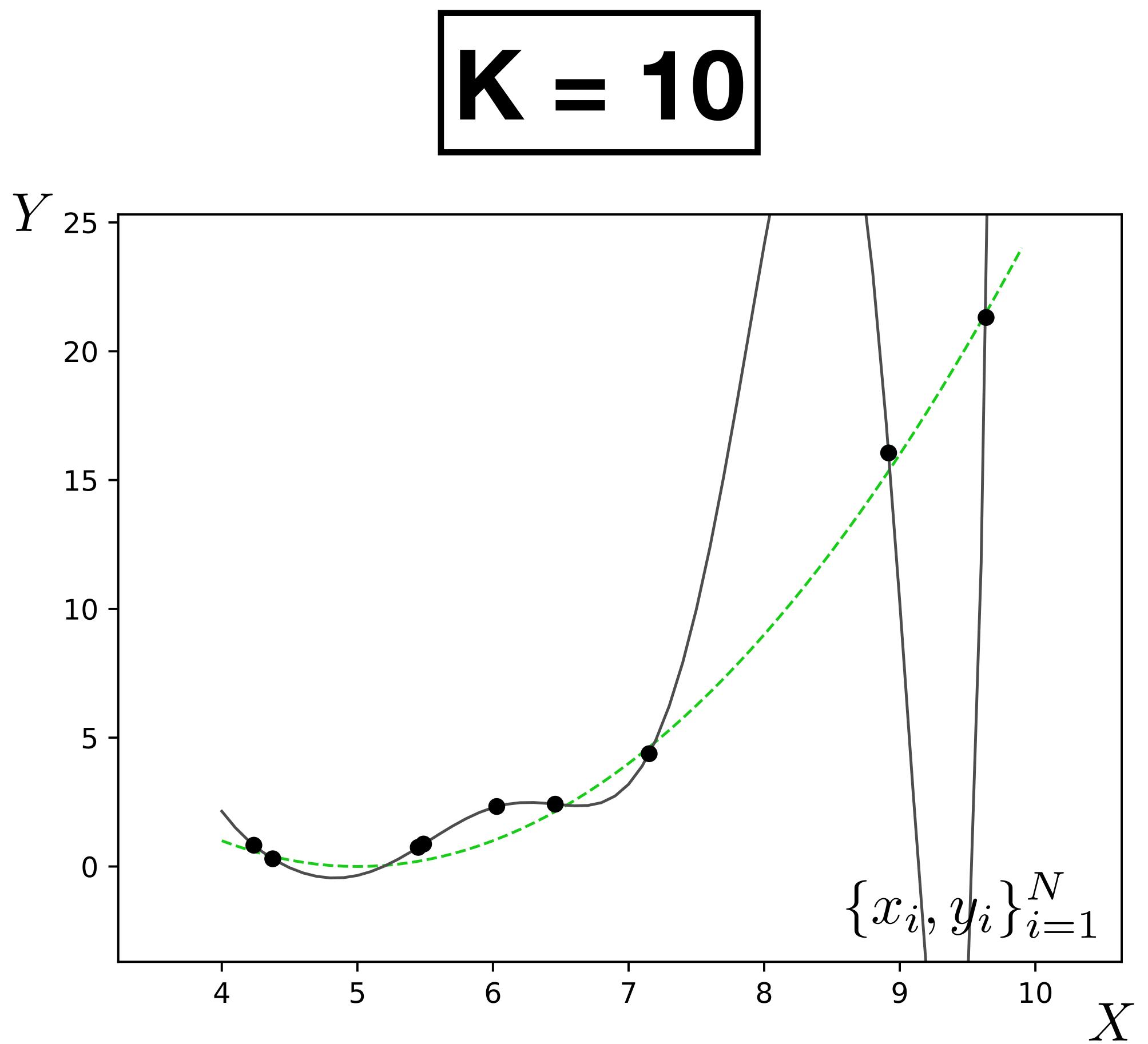
$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

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?

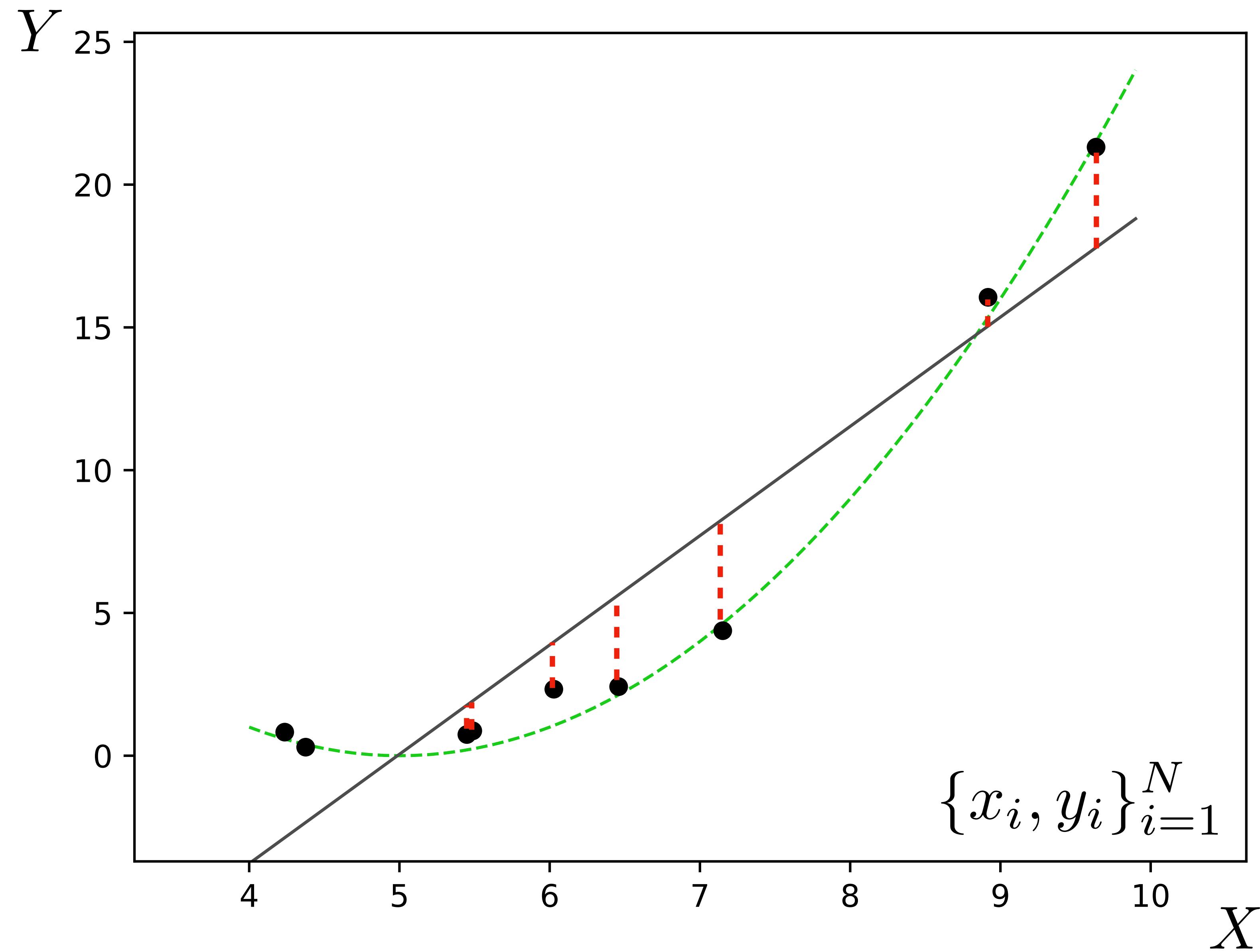


$$f_{\theta}(x) = \sum_{k=0}^{K} \theta_k x^k$$

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.

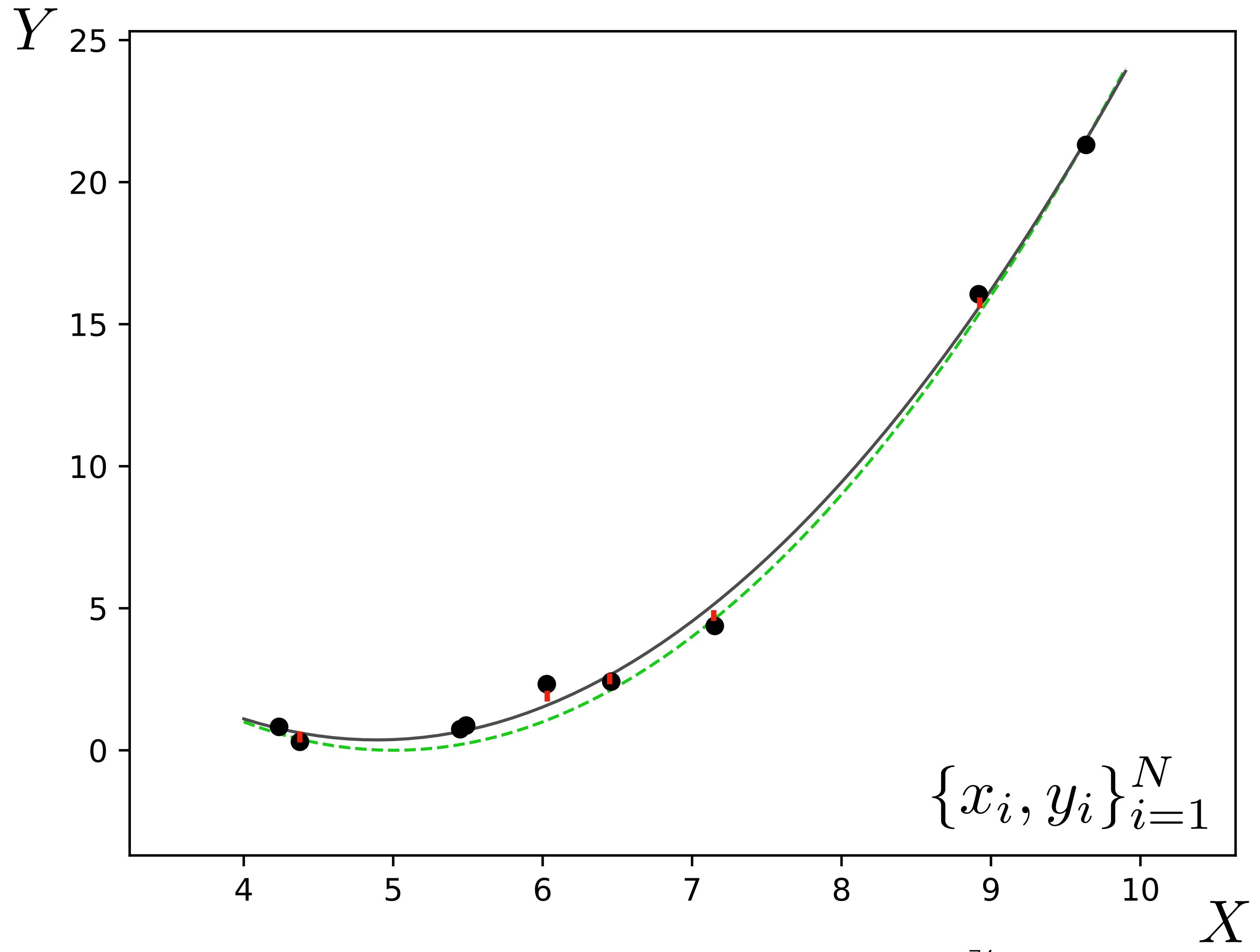
K = 1



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

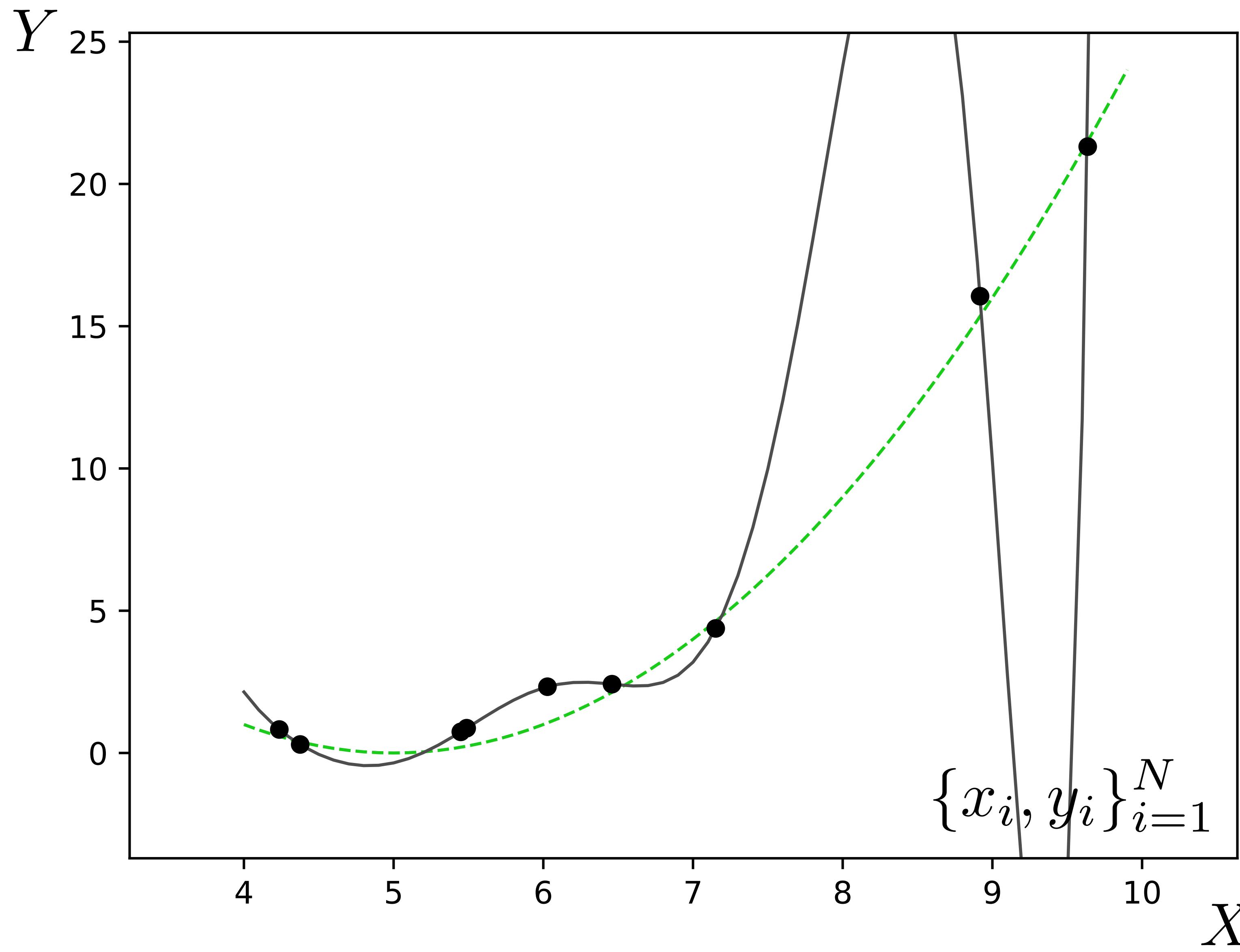
An underfit model will have high **error** on the training points. This error is known as **approximation error**.

K = 2



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

K = 10

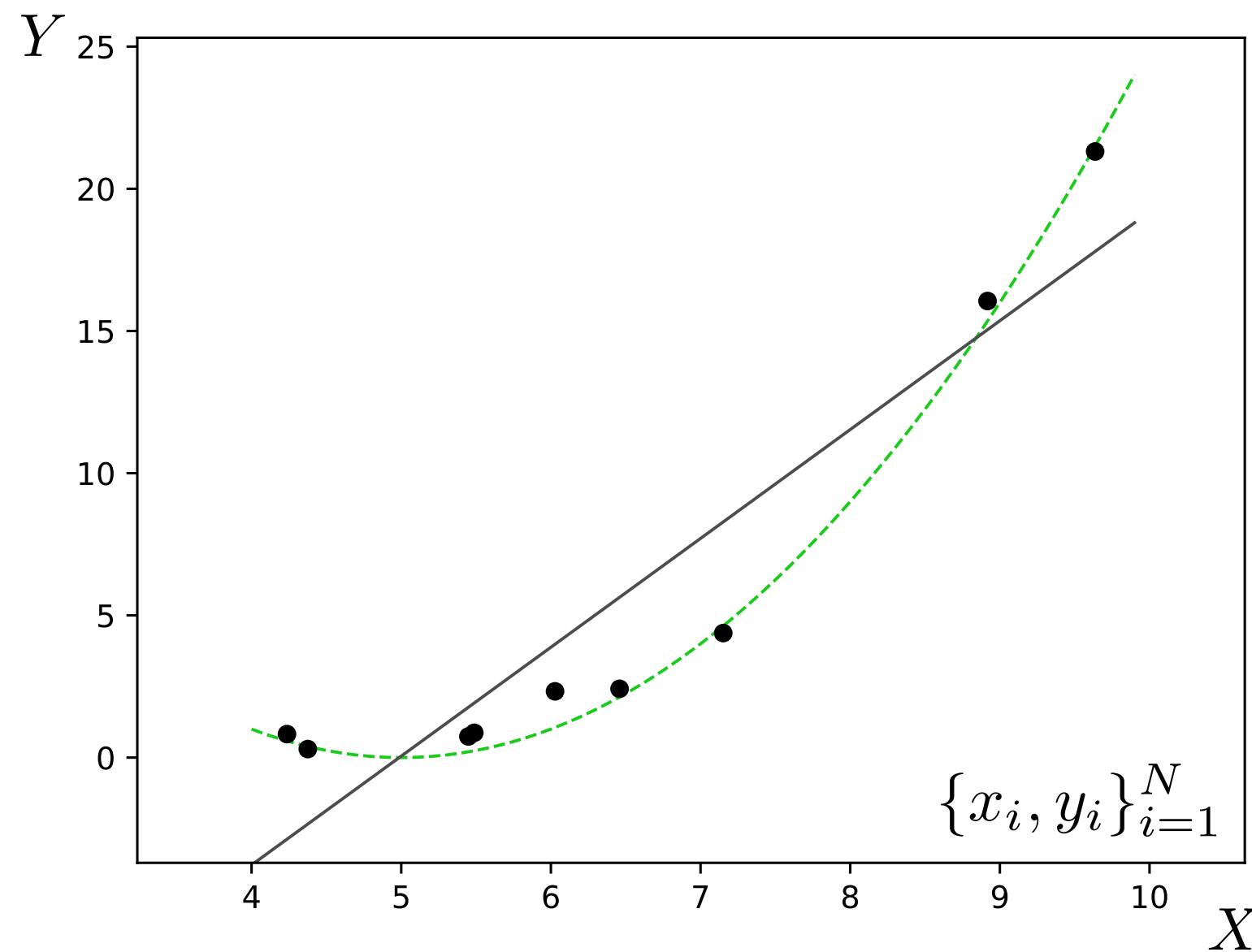


Now we have zero approximation error – the curve passes exactly through each training point.

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).

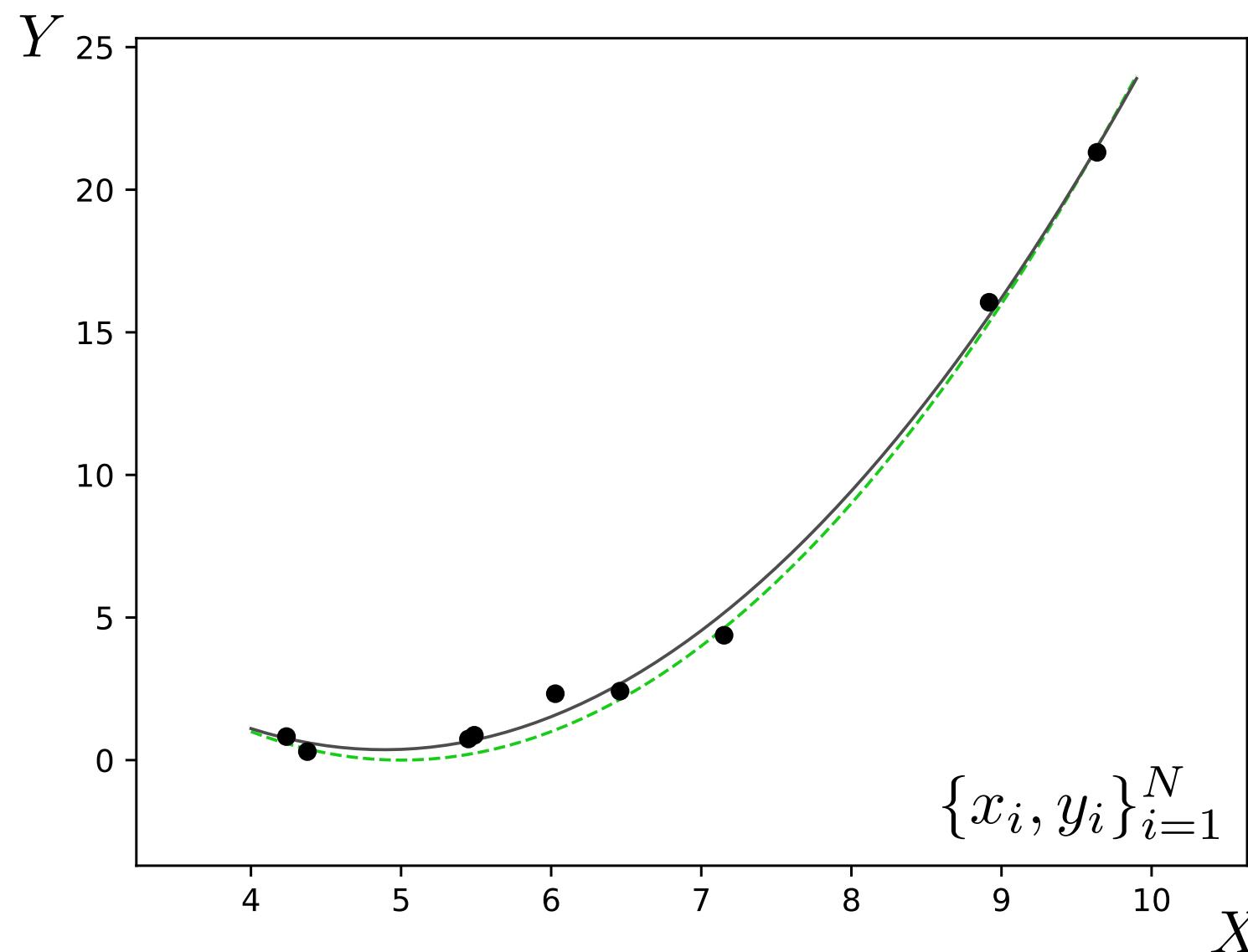
Underfitting

$$K = 1$$



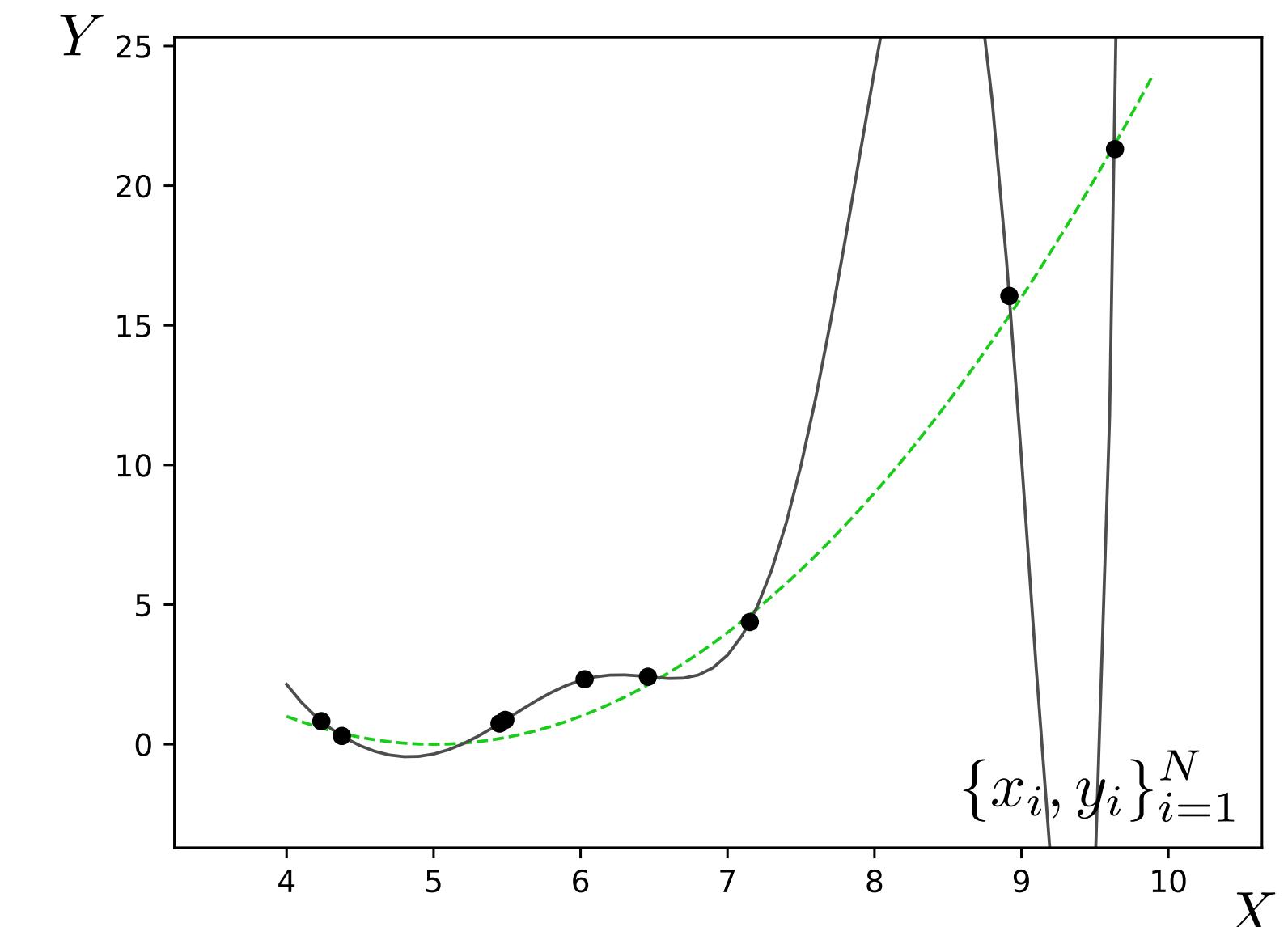
Appropriate model

$$K = 2$$



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

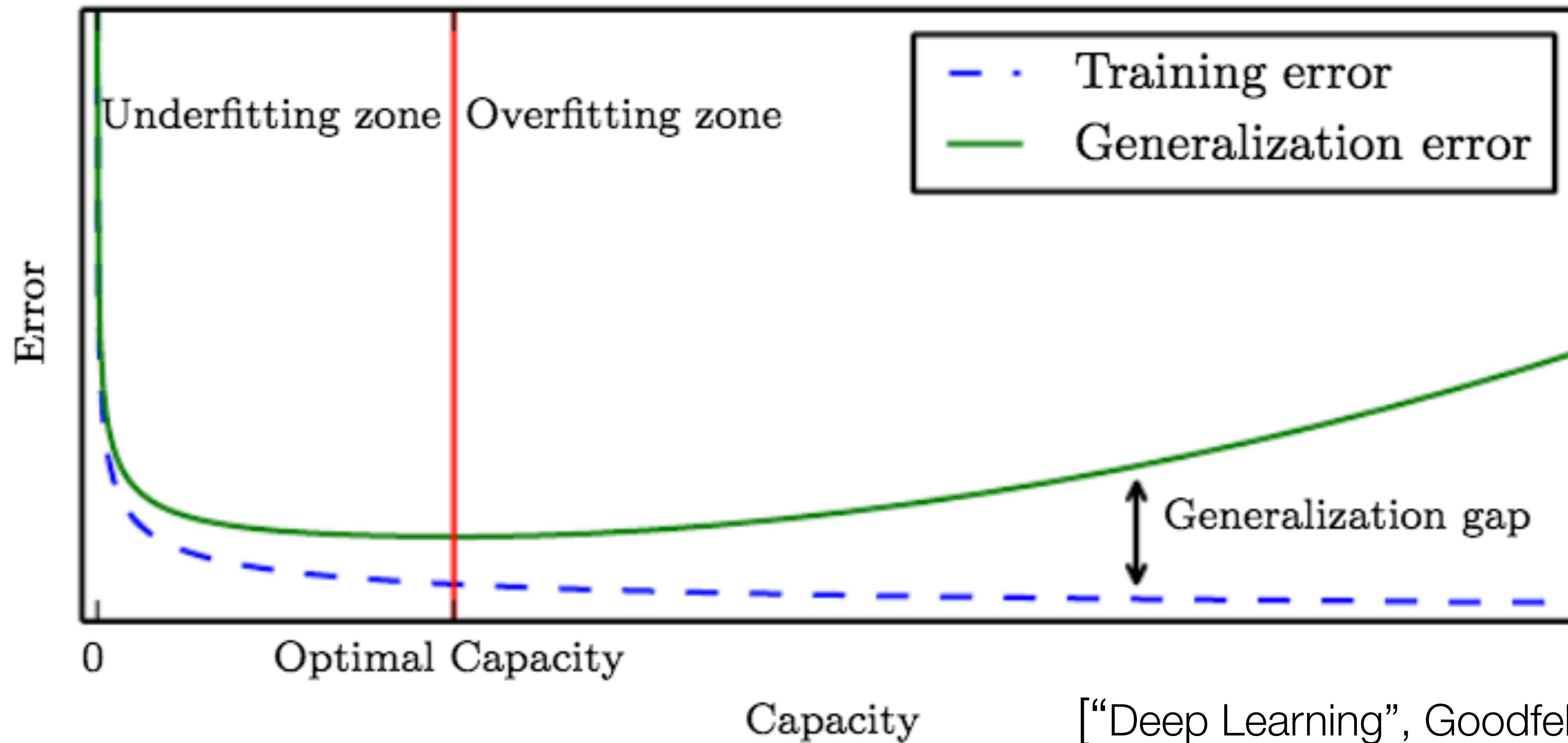
We need to control the **capacity** of the model (e.g., use the appropriate number of free parameters).

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.

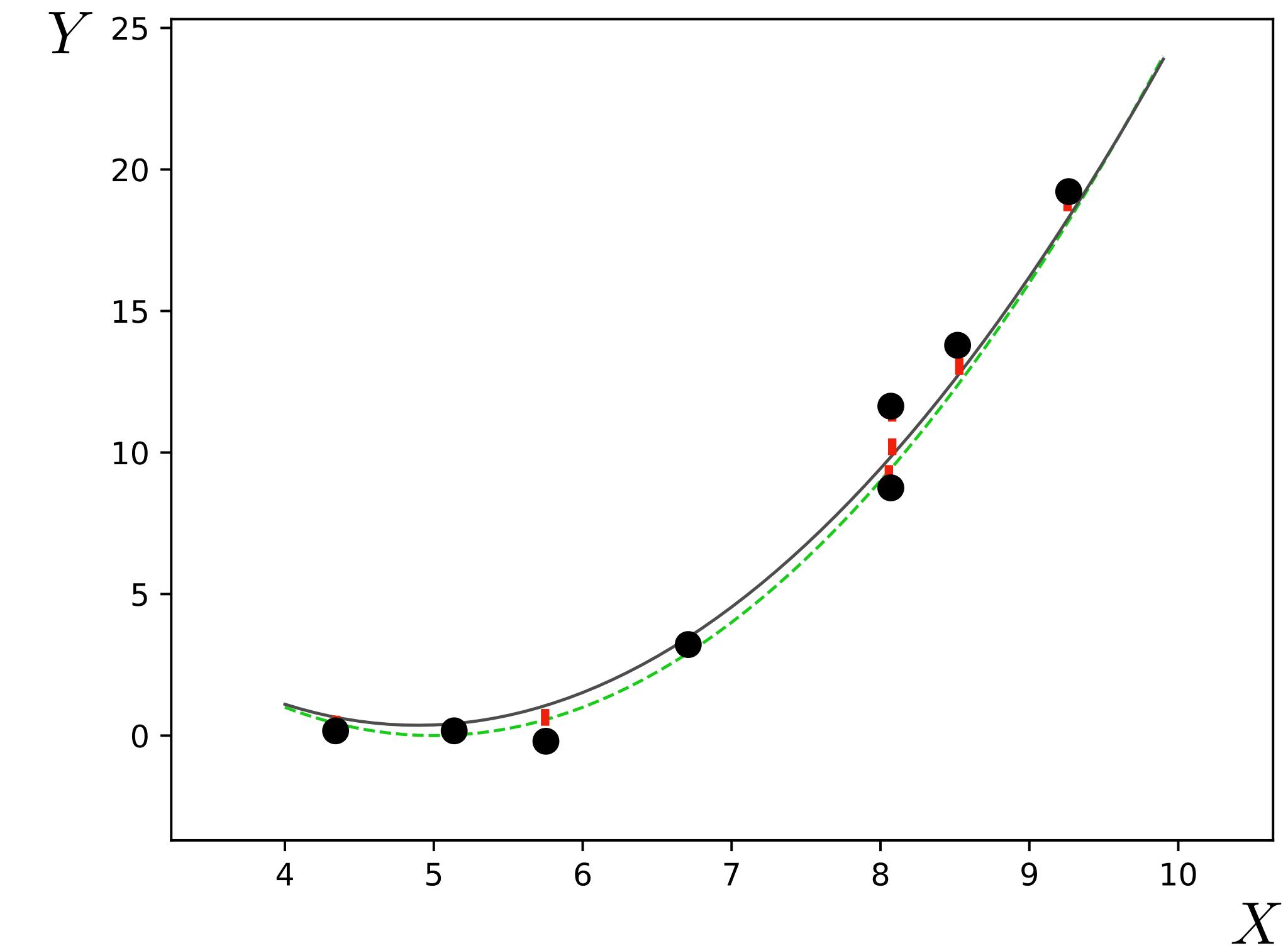
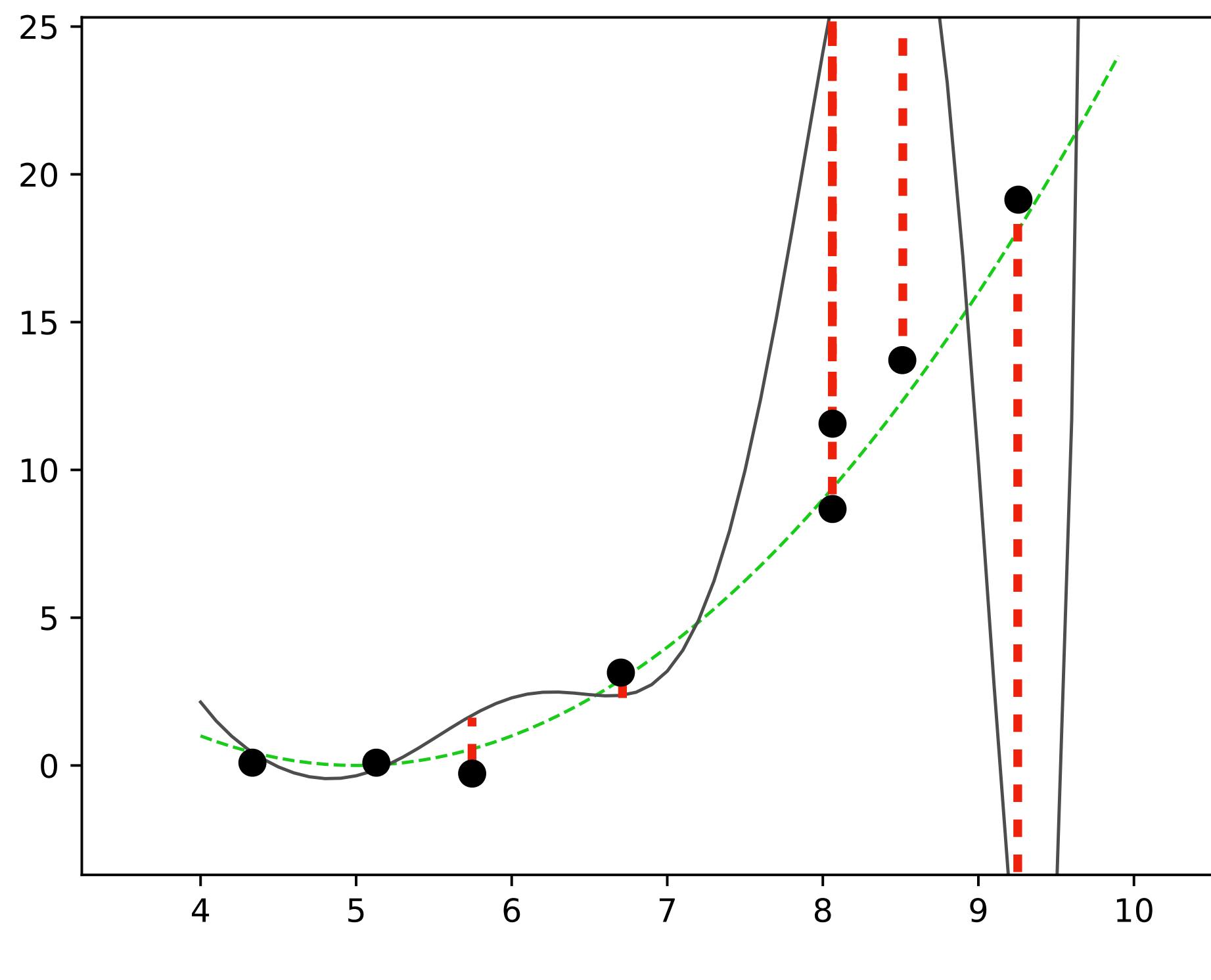
Simple models have low capacity.

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 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:

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

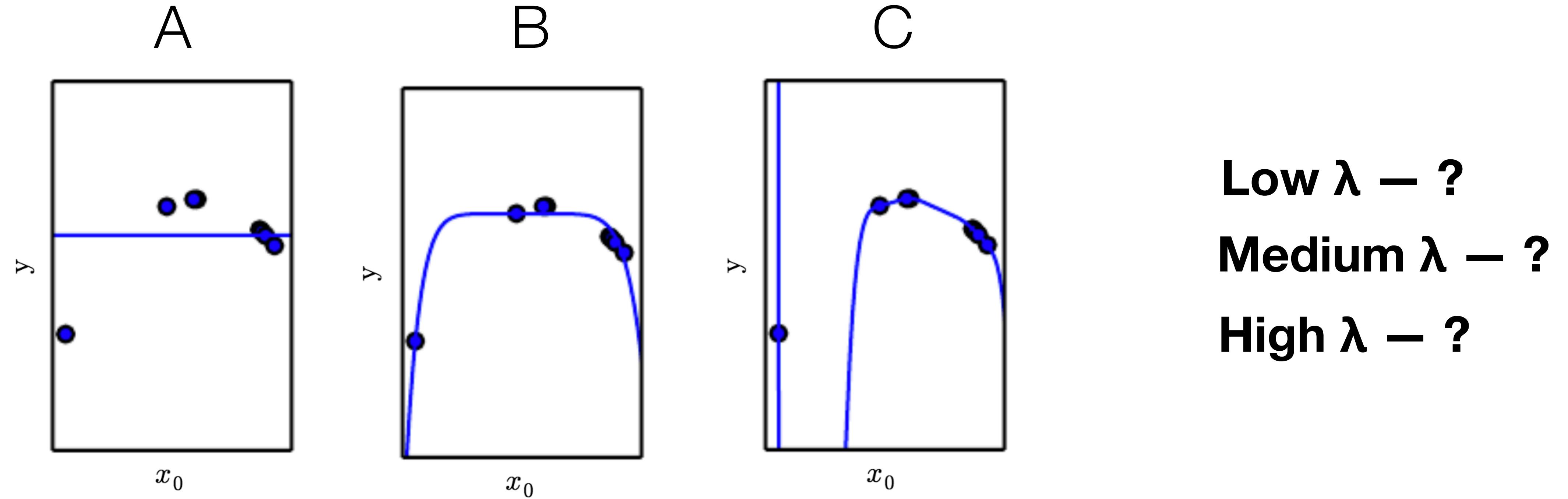
Regularized least squares

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

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

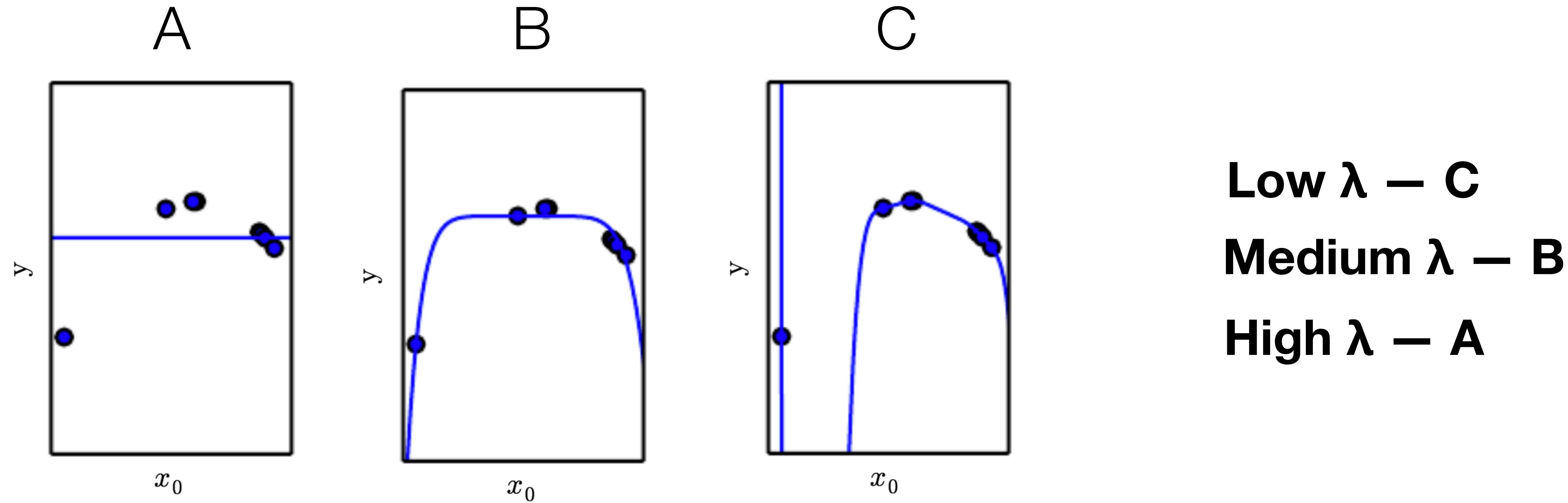
ridge regression, a.k.a., **Tikhonov regularization**

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



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

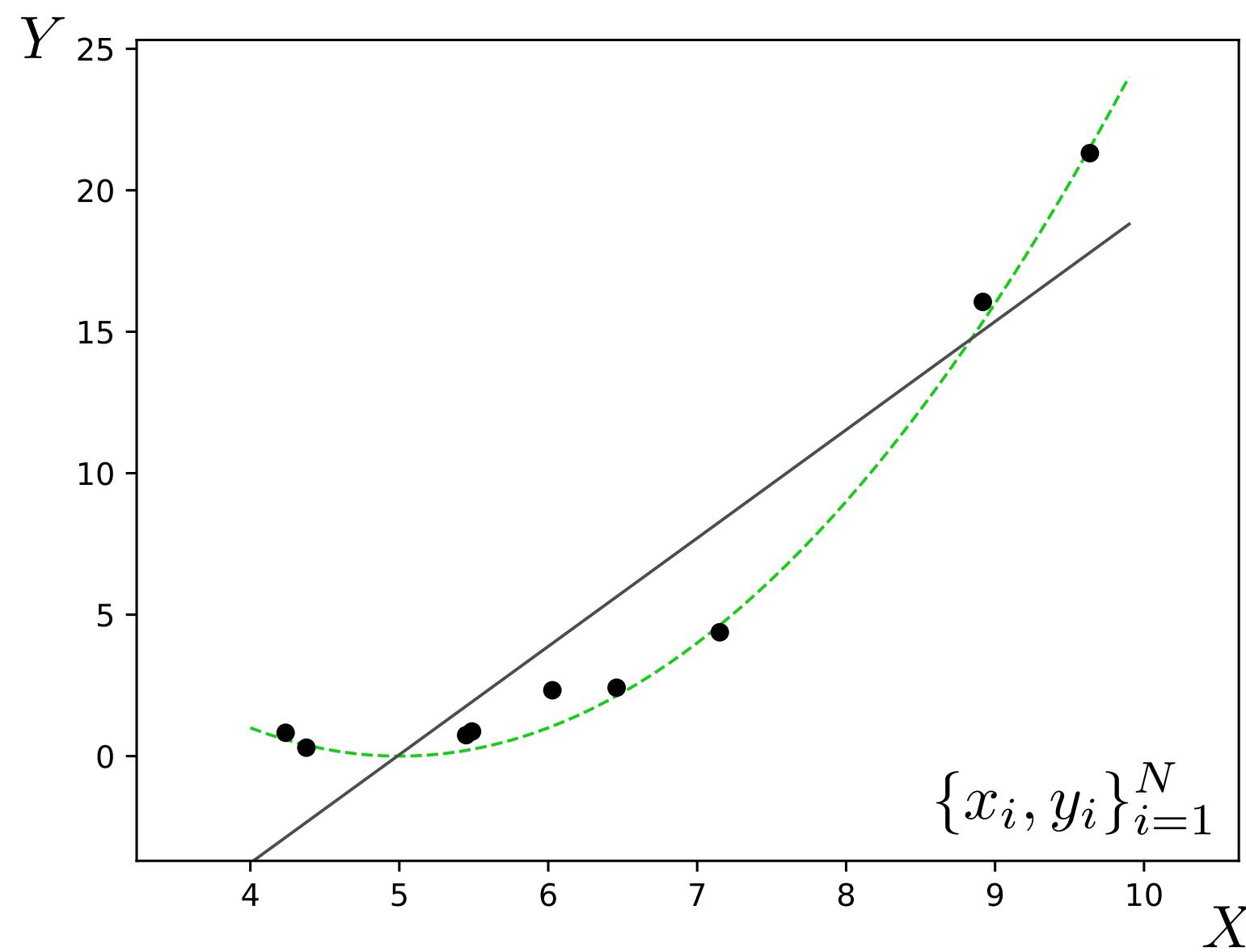
$$\theta^* = \arg \min_{\theta} \sum_{I=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + \lambda \|\theta\|_2^2$$



[Adapted from “Deep Learning”, Goodfellow et al.]
84

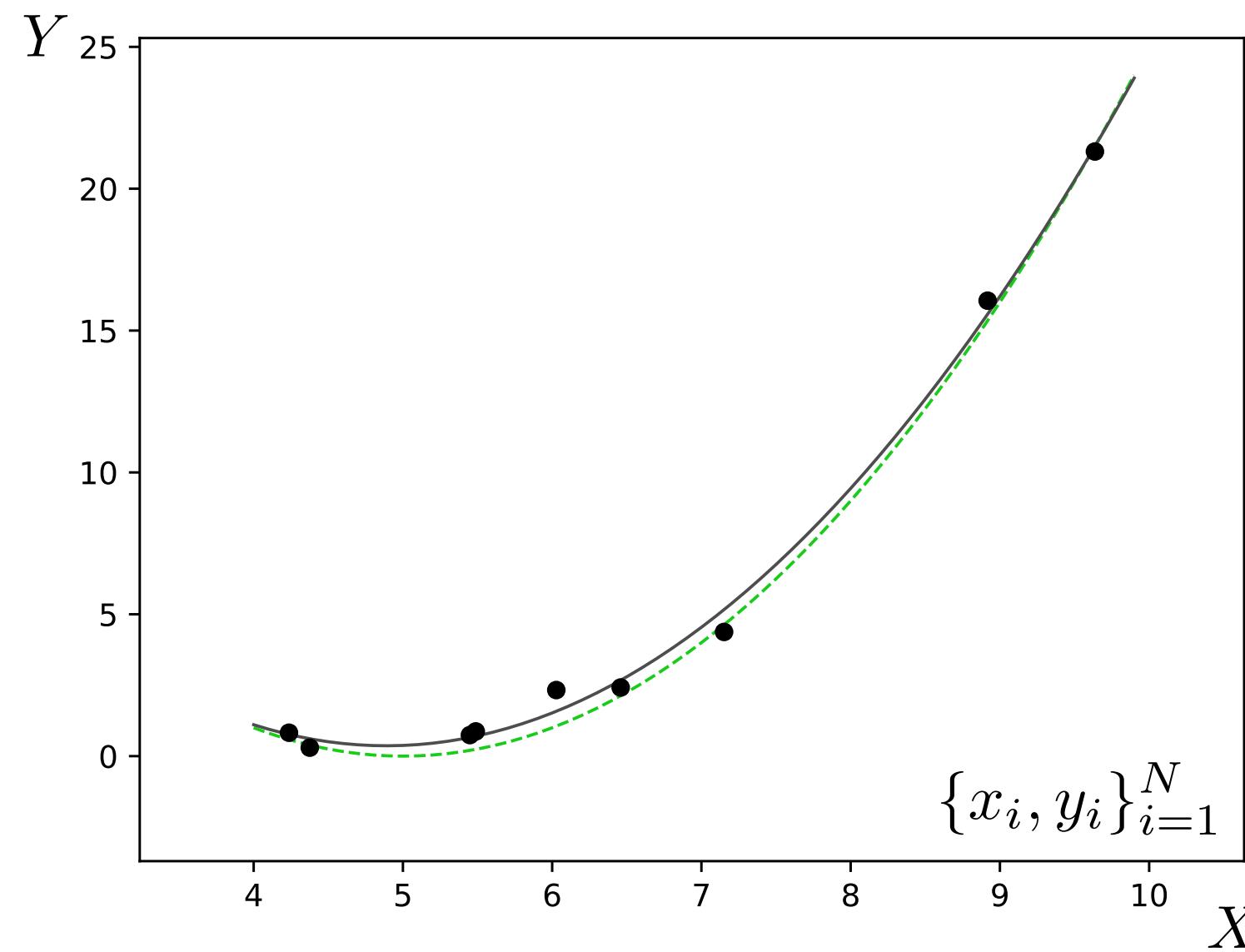
Underfitting

$$K = 1$$



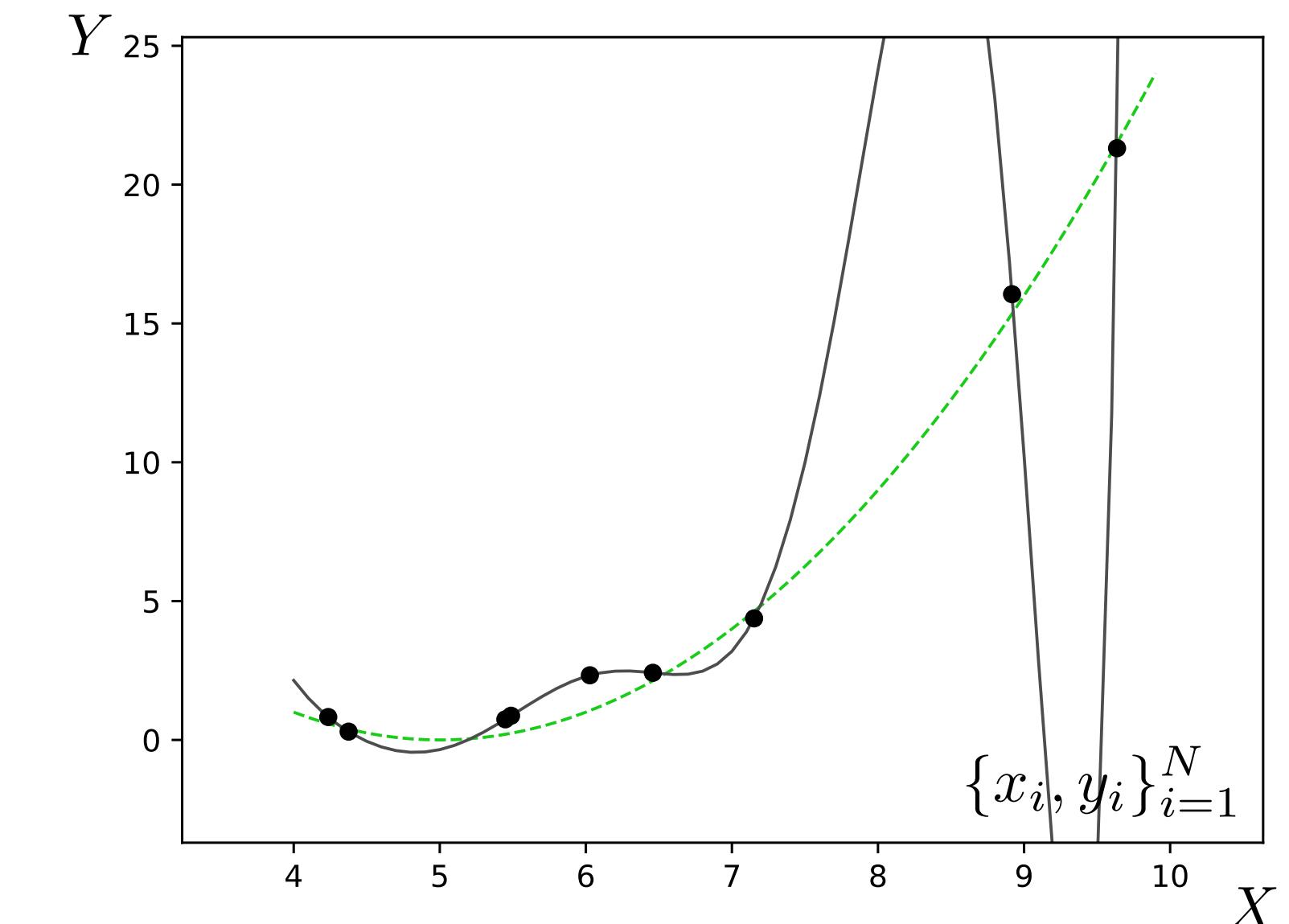
Appropriate model

$$K = 2$$



Overfitting

$$K = 10$$



Simple model

Doesn't fit the training data

Simple model

Fits the training data

Complex model

Fits the training data

Image classification

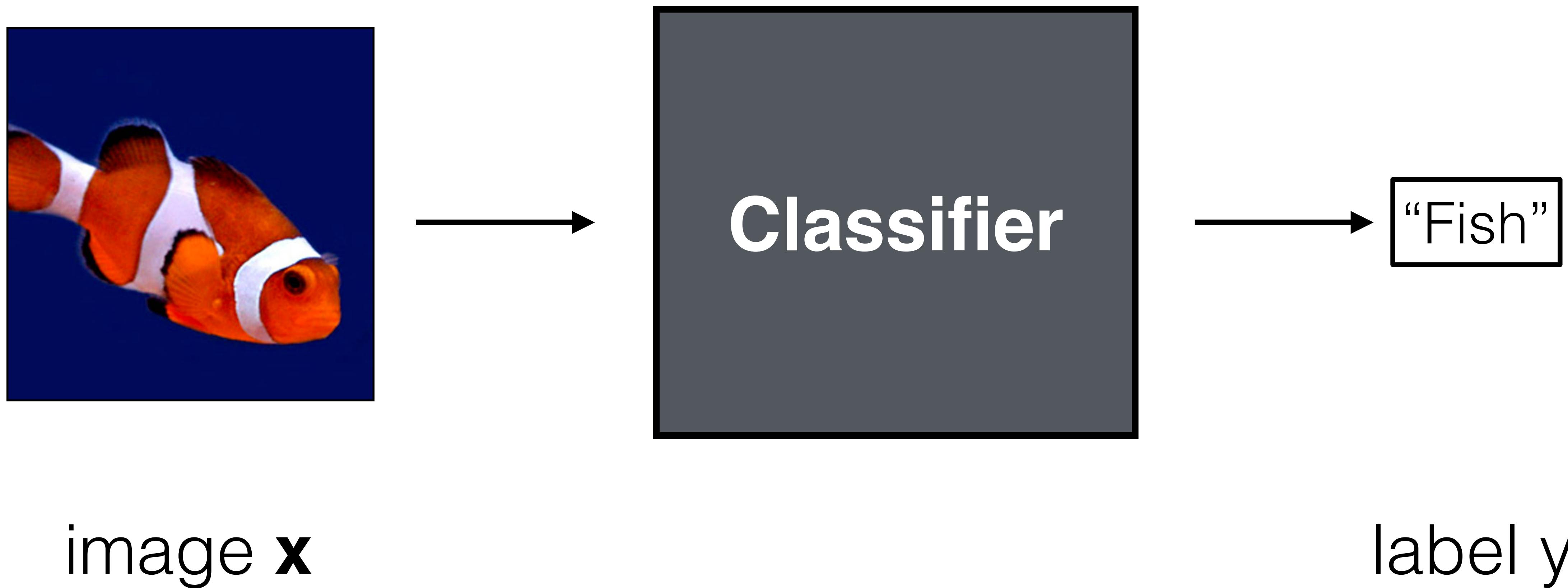


image x

label y

Image classification



“Fish”

image x

label y

Image classification

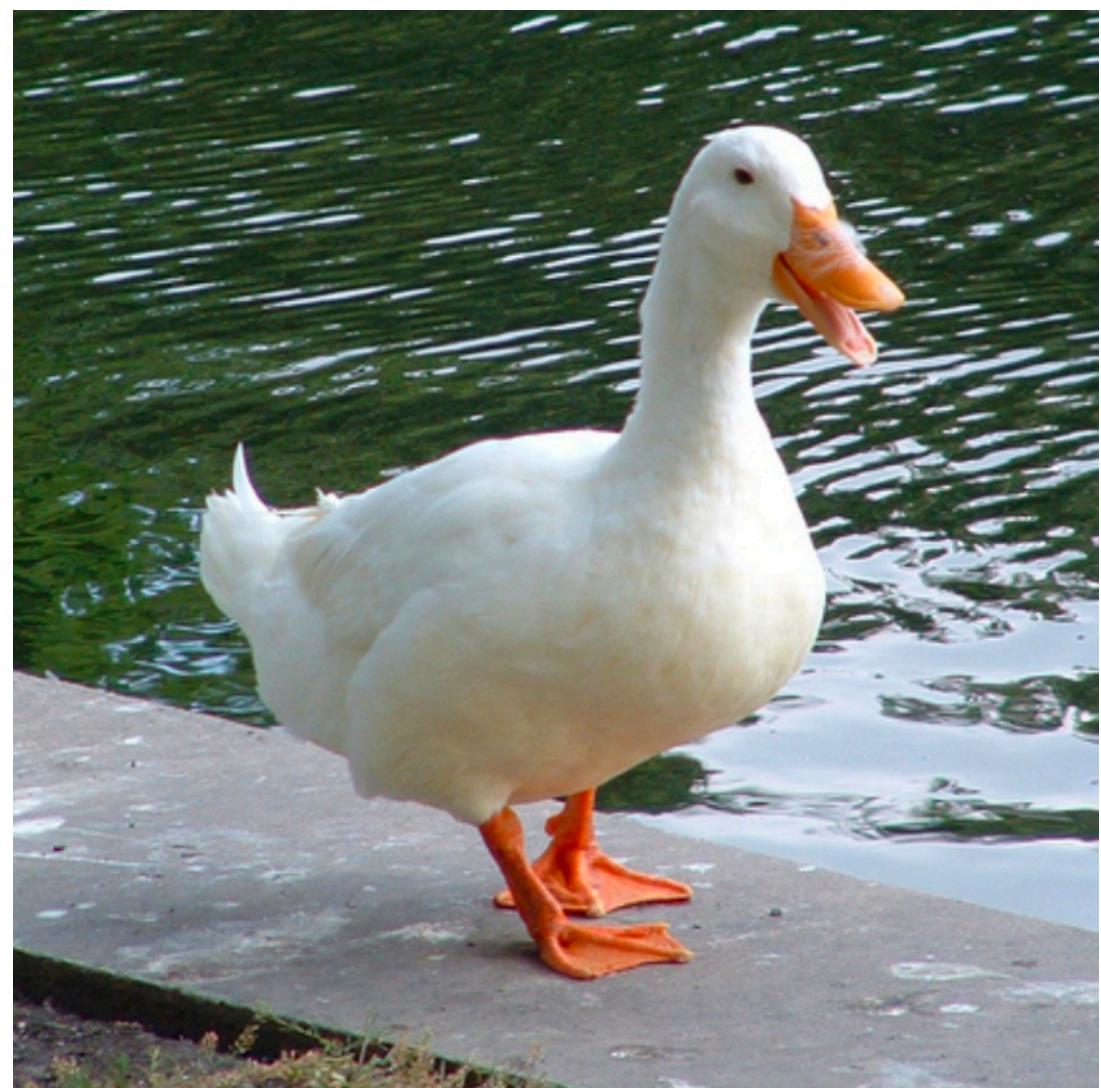


“Fish”

image x

label y

Image classification



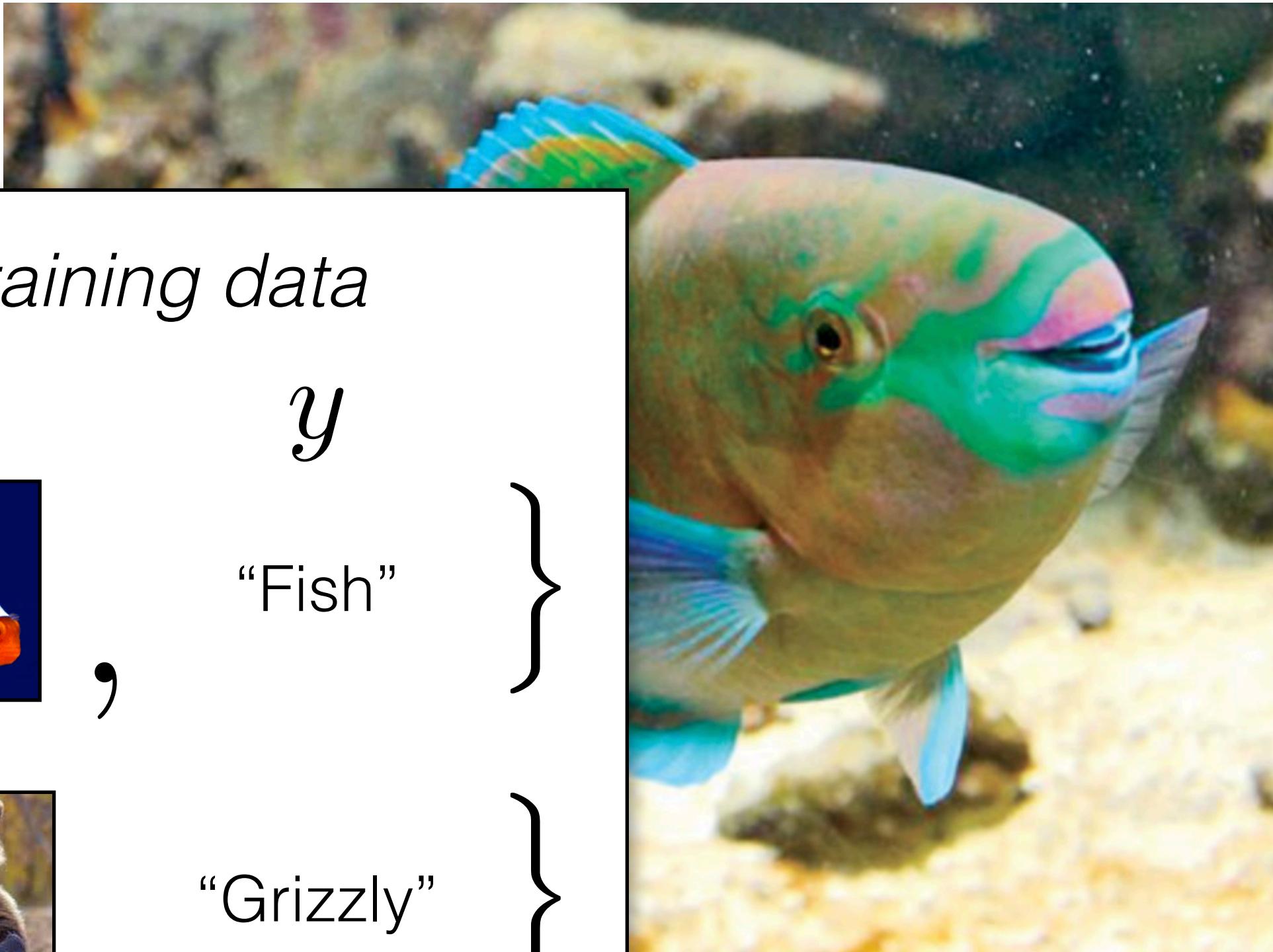
“Duck”

:

image **x**

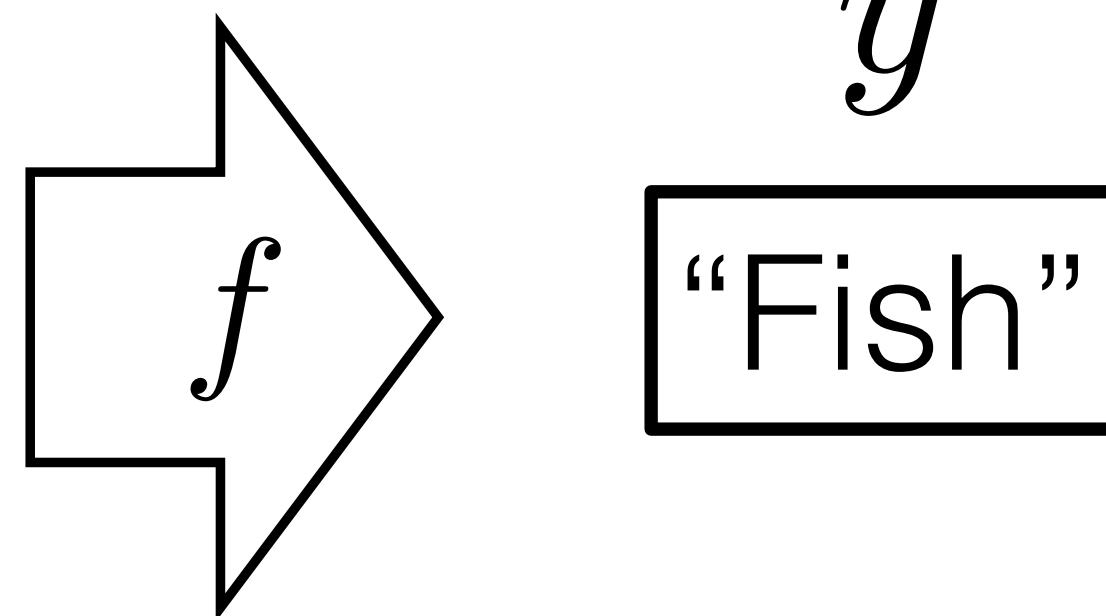
label **y**

X



Training data

x	y
{  , }	“Fish”
{  , }	“Grizzly”
{  , }	“Chameleon”
:	

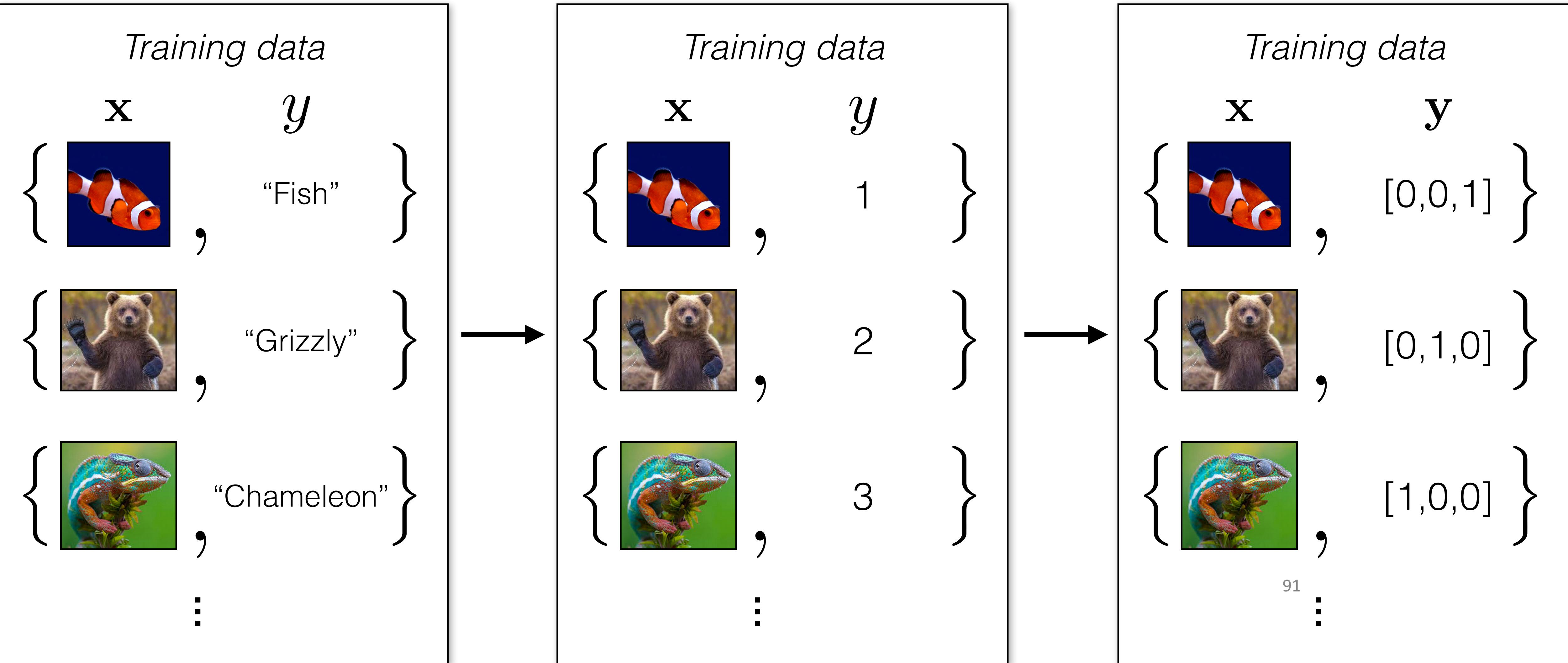


$$\arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i), y_i)$$

90

How to represent class labels?

One-hot vector



What should the loss be?

0-1 loss (number of misclassifications)

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

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

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{k=1}^K (y_k - \hat{y}_k)^2 \quad \leftarrow \text{easy to optimize, but crude approximation}$$

Cross entropy

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad \leftarrow \text{easy to optimize, good approximation}$$

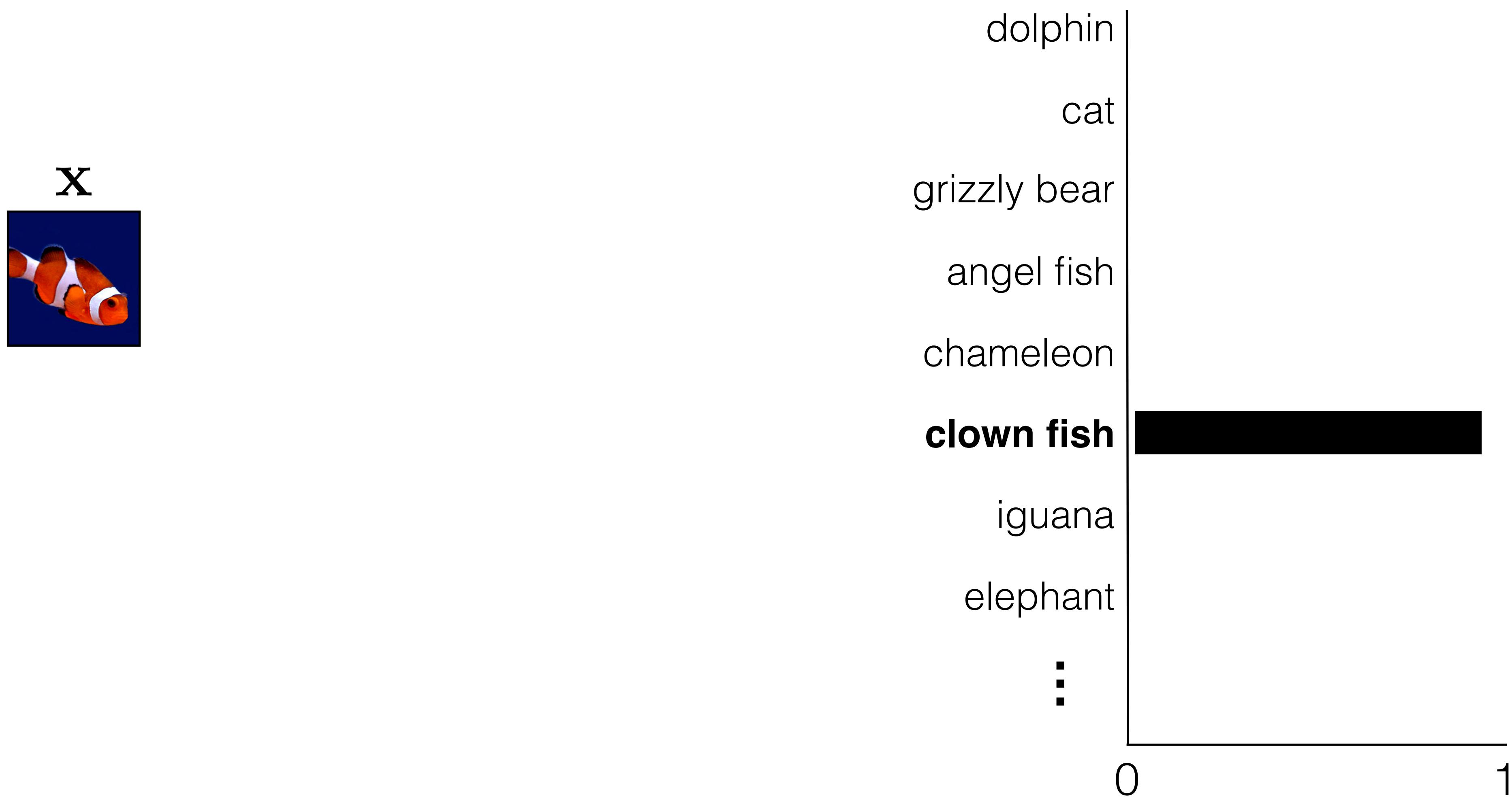
Ground truth label y

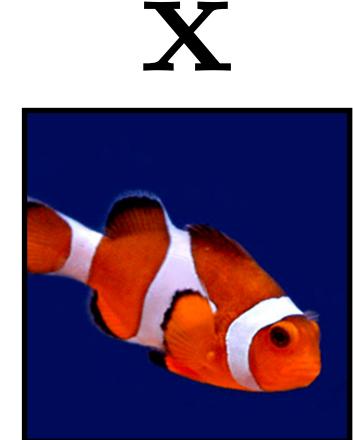
x



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

Ground truth label **y**





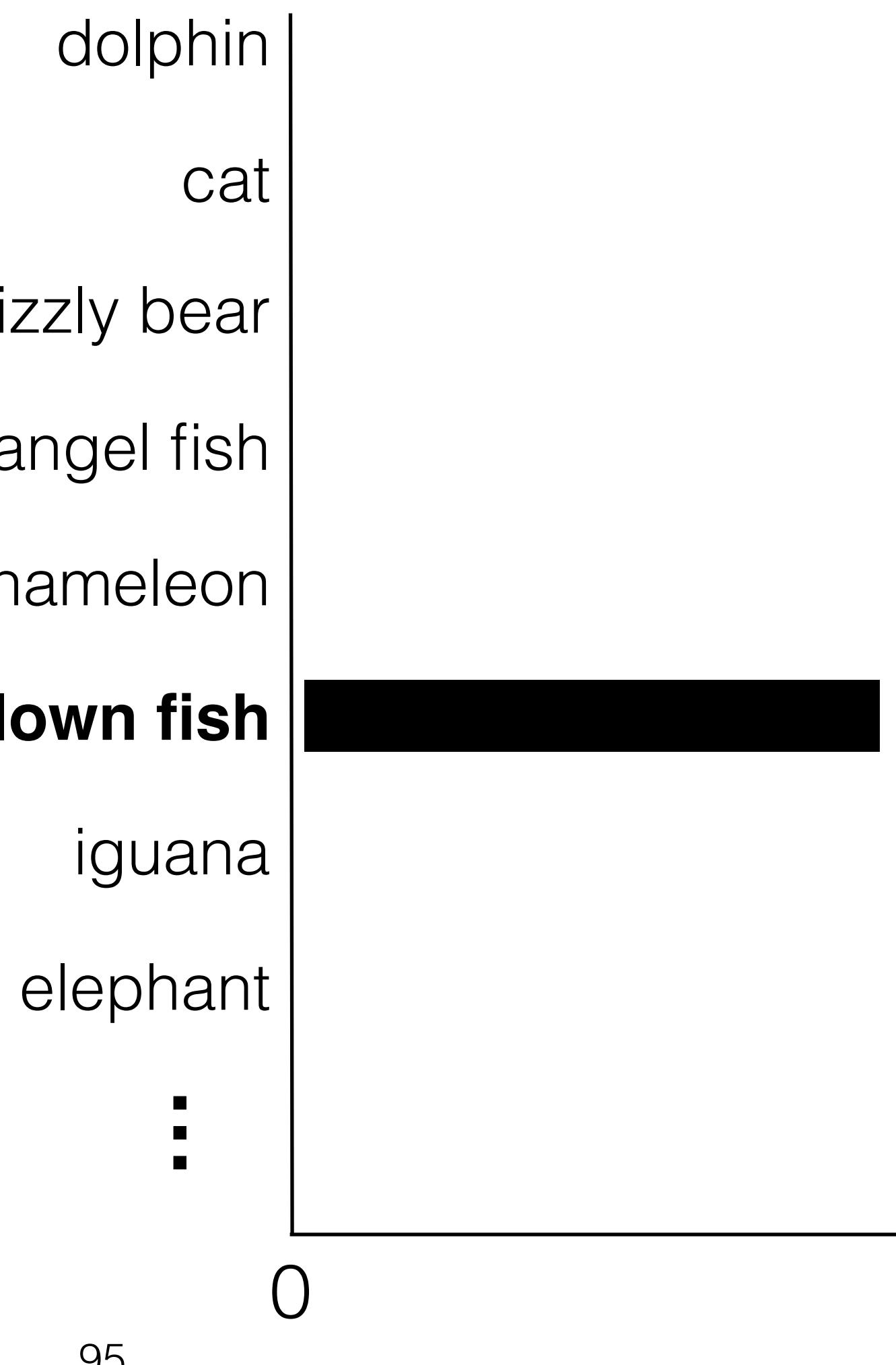
X

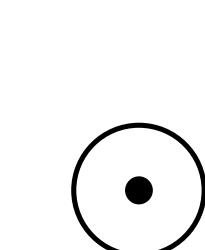
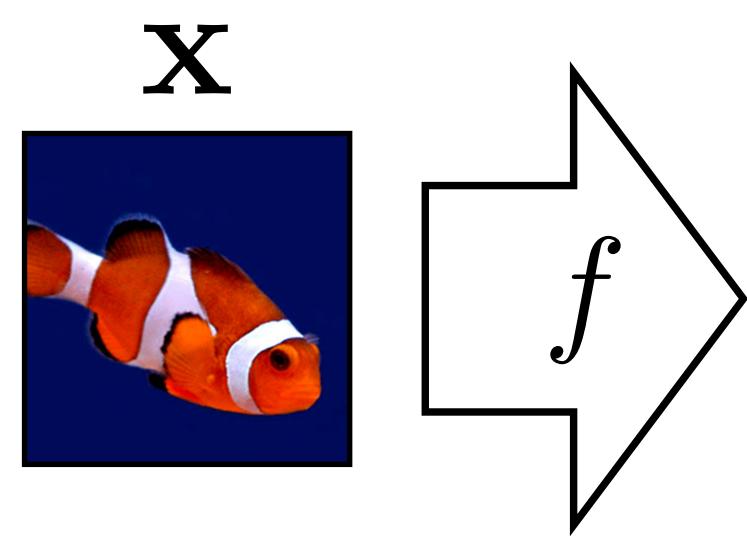
Prediction \hat{y}

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



Ground truth label y





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

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

dolphin	█
cat	█
grizzly bear	█
angel fish	█
chameleon	█
clown fish	██████████
iguana	█
elephant	█
⋮	⋮

0

1

96

Ground truth label \mathbf{y}

dolphin	
cat	
grizzly bear	
angel fish	
chameleon	
clown fish	██████████
iguana	
elephant	
⋮	⋮

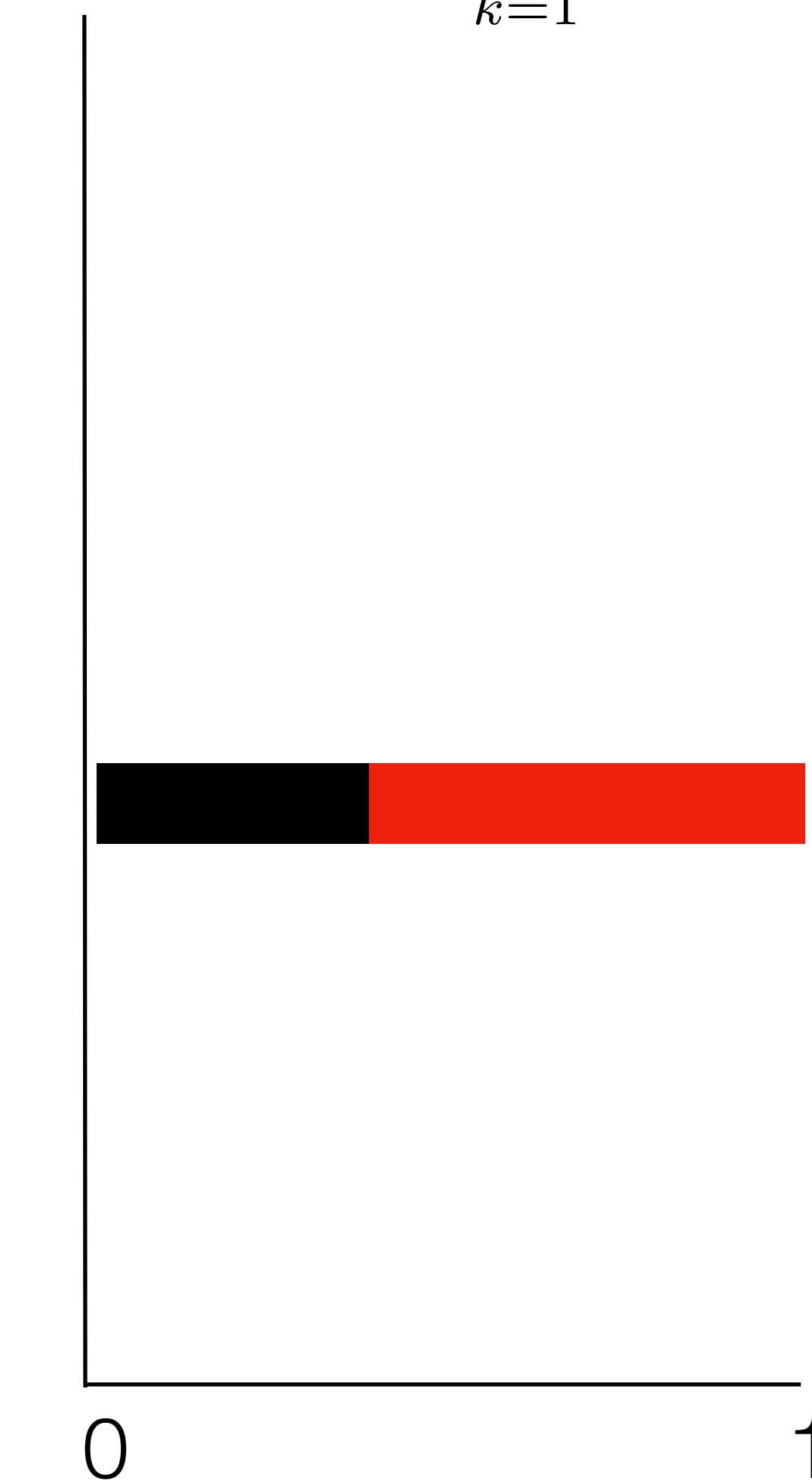
0

1

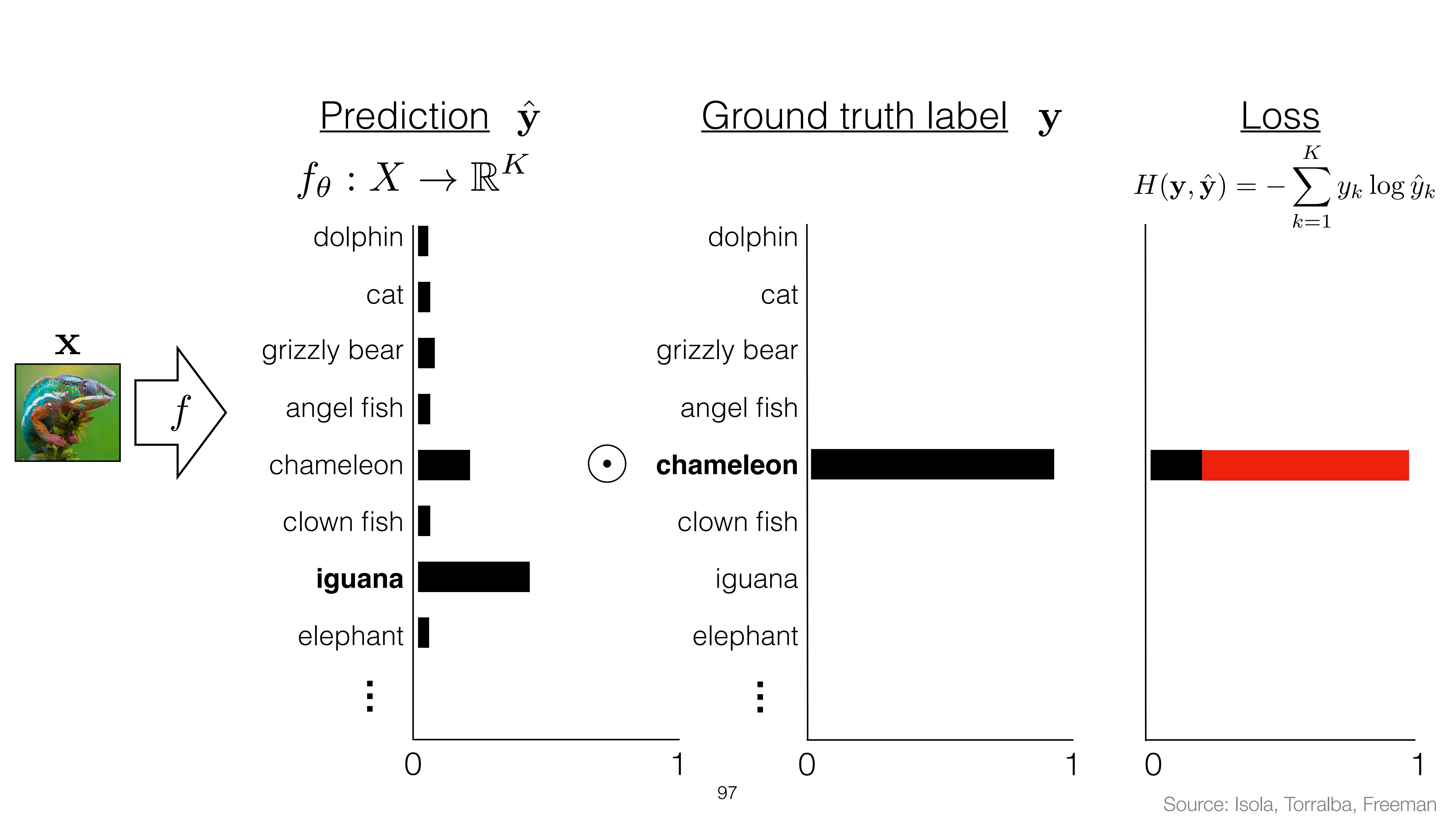
Source: Isola, Torralba, Freeman

Loss

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



1



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

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

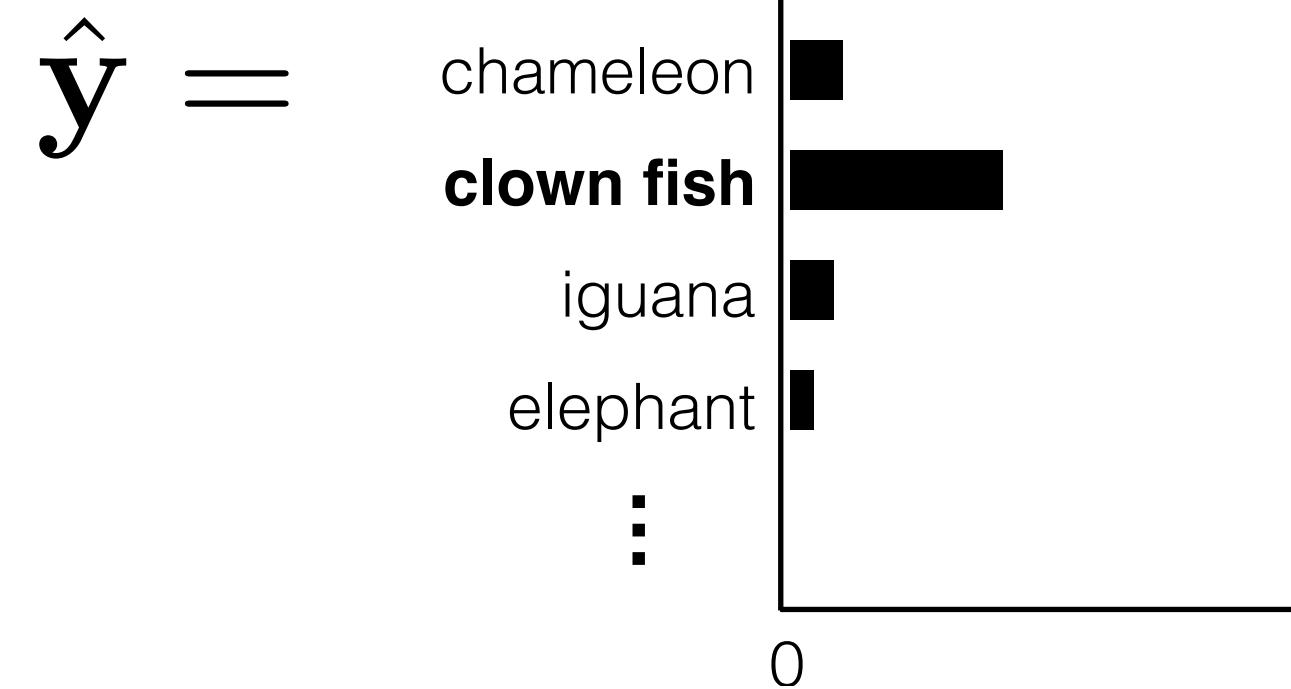
$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

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

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

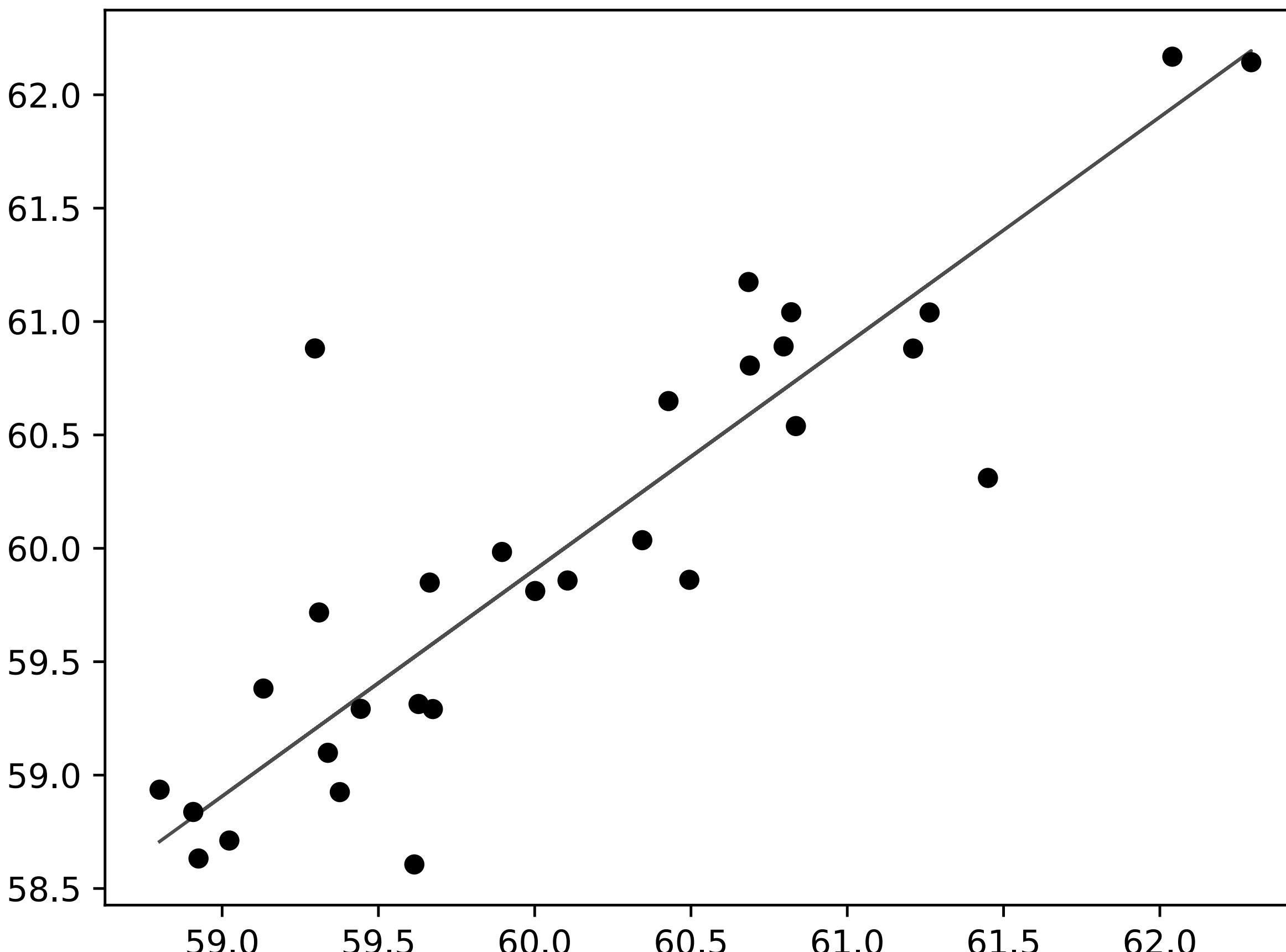
← squash into a non-negative vector that sums to 1
— i.e. **a probability density function**

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



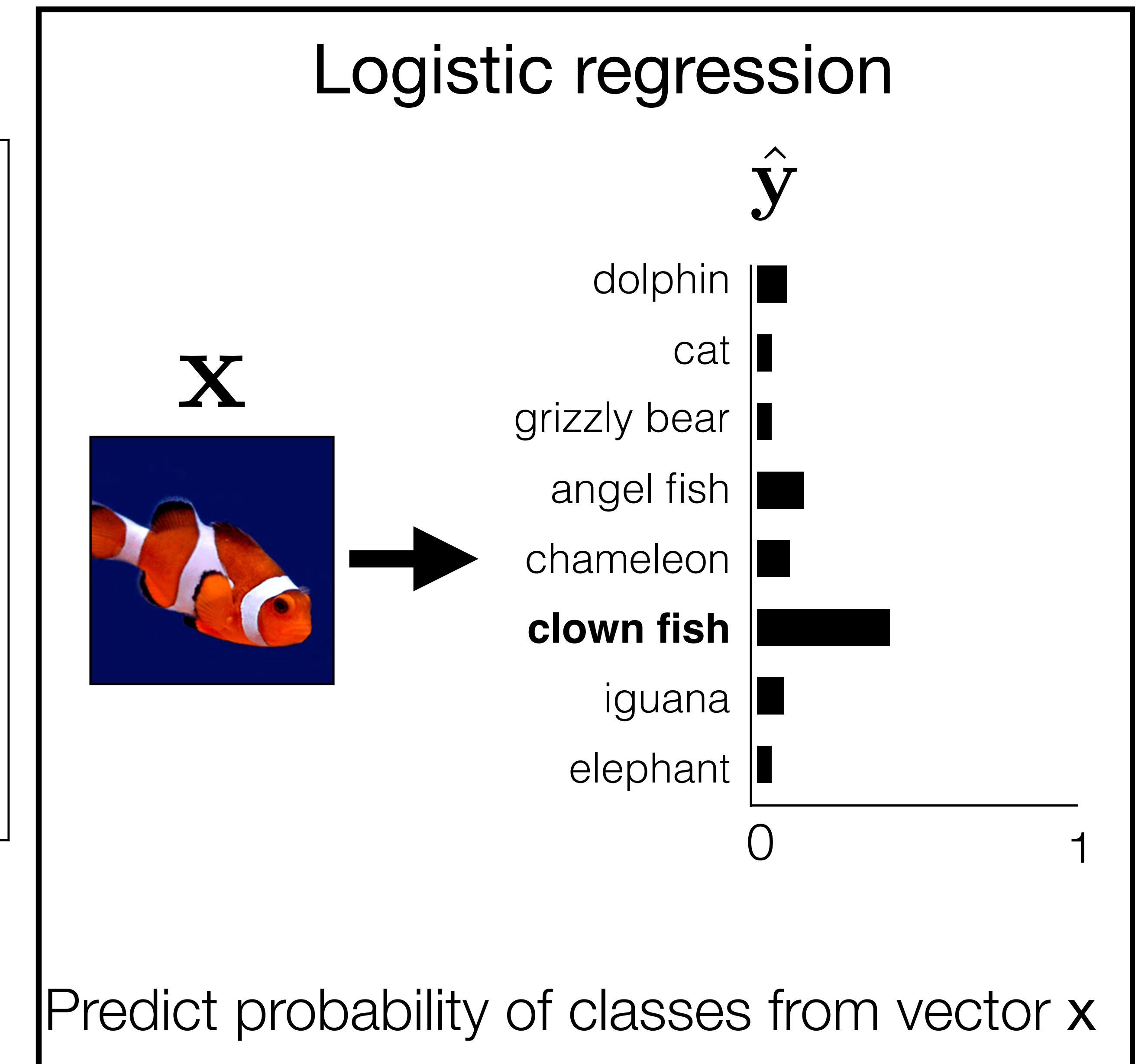
98

Linear regression



Predict scalar y from vector x

Logistic regression

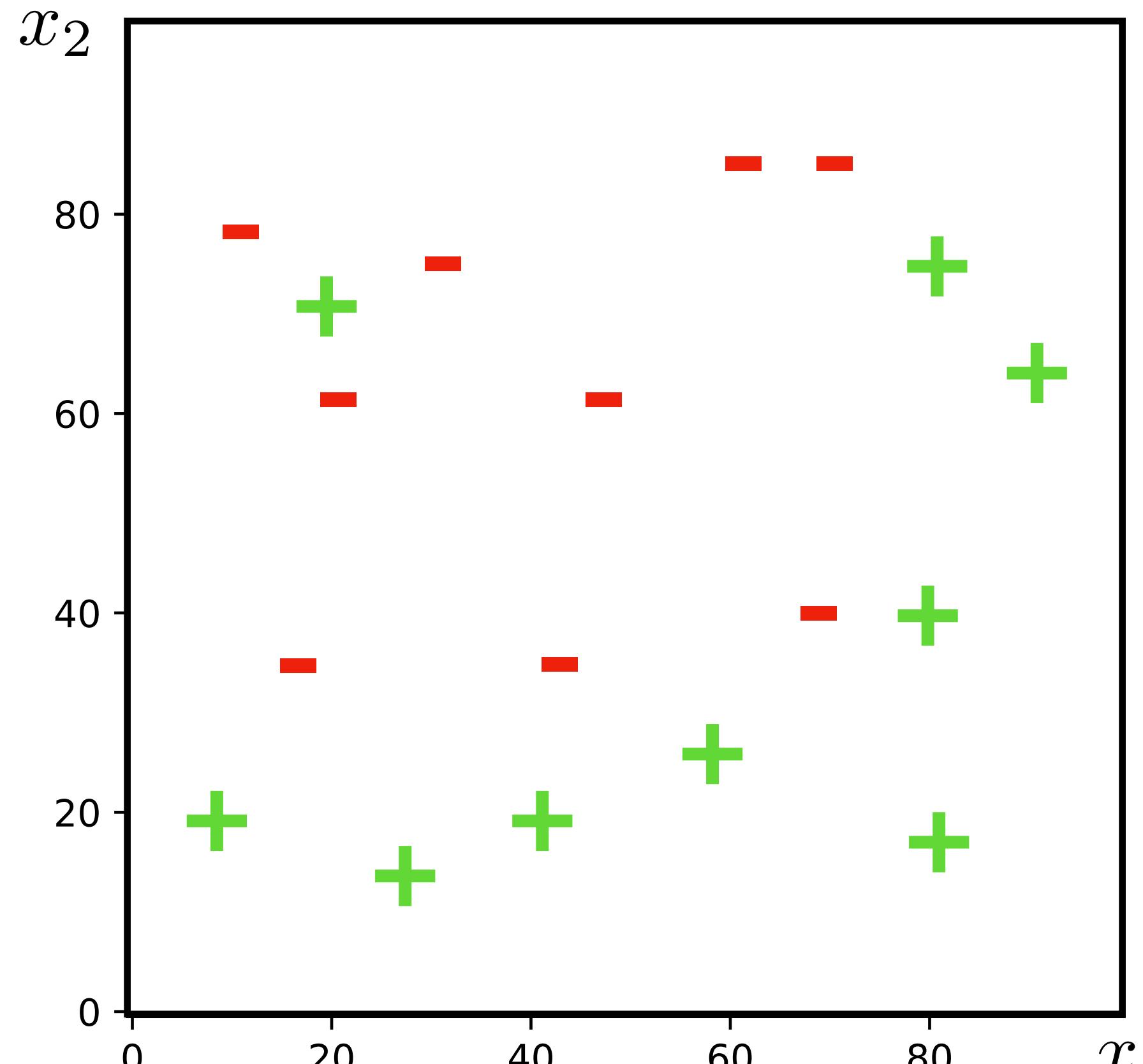


Predict probability of classes from vector x

Putting these ideas together: linear classifiers

Linear decision boundaries

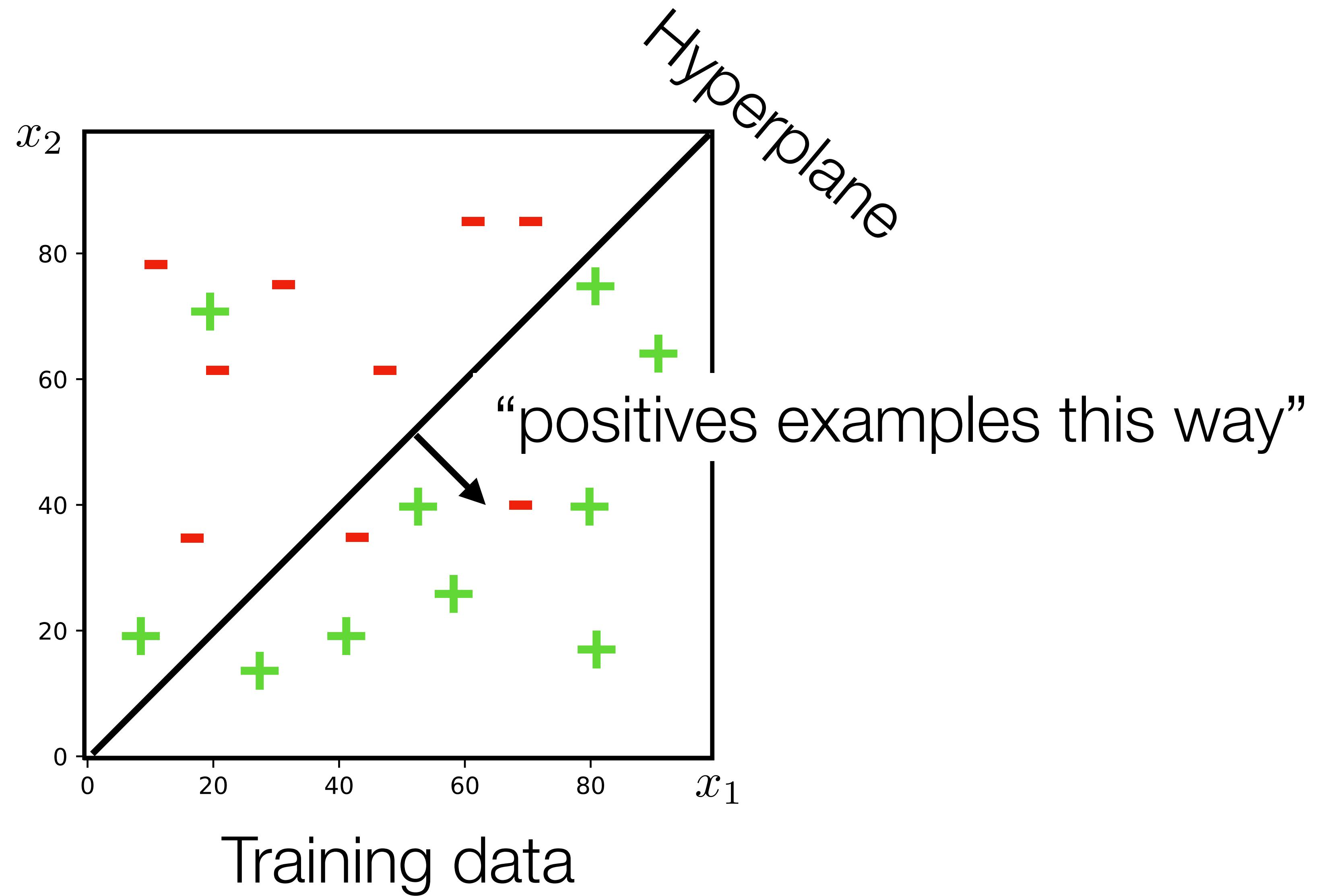
Consider a **binary** classification problem.



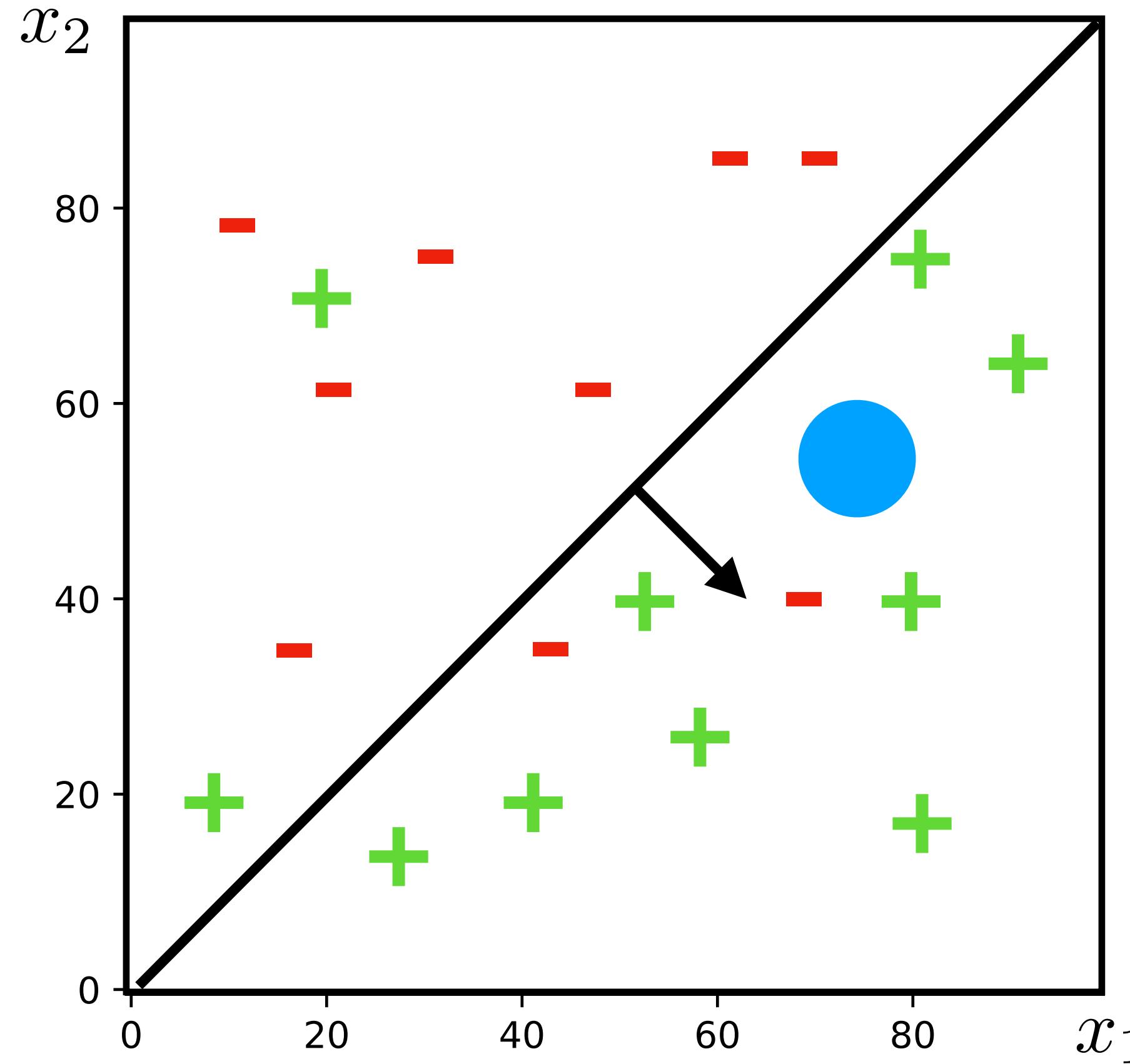
Training data

Linear decision boundaries

Consider a **binary** classification problem.



Linear decision boundaries



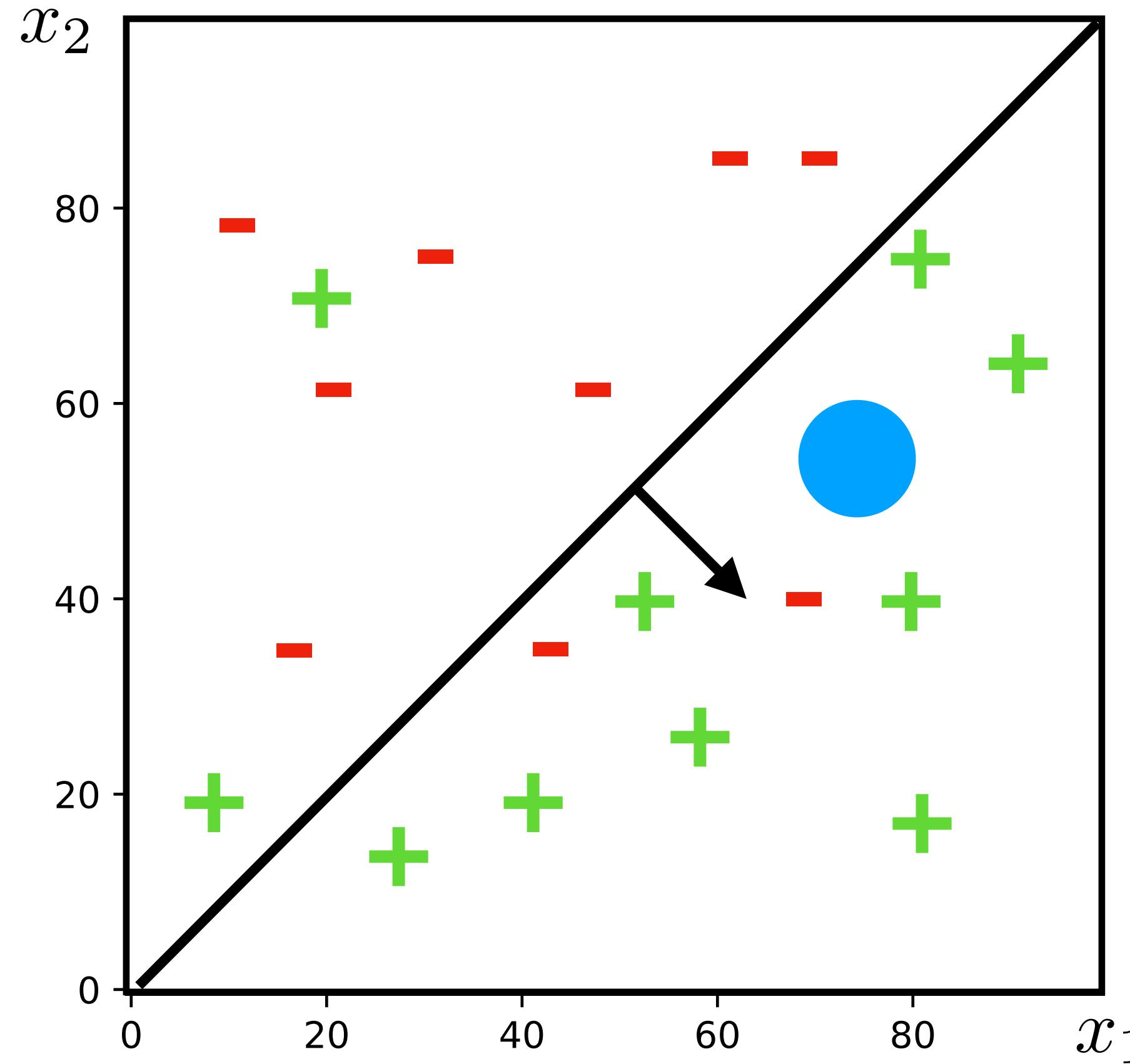
“What label is this point?”

Which side of the hyperplane is \mathbf{x} on?

$$\hat{y} = \mathbf{x}^T \mathbf{w} + b$$

$$g(\hat{y}) = \begin{cases} 1, & \text{if } \hat{y} > 0 \\ 0, & \text{otherwise} \end{cases}$$

Linear decision boundaries



Notational simplification:

Can get rid of b by “tacking on” a 1 to x

$$\hat{y} = \tilde{\mathbf{x}}^\top \tilde{\mathbf{w}}$$

$$\tilde{\mathbf{x}} = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$\tilde{\mathbf{w}} = \begin{bmatrix} w \\ b \end{bmatrix}$$

“What label is point?”

Multiway classification

For a k -class problem, we'll make a matrix “stacking” k hyperplanes as rows of a matrix $W \in \mathbb{R}^{k \times d}$.

$$W = \begin{bmatrix} \text{---} & w_1 & \text{---} \\ \text{---} & w_2 & \text{---} \\ \text{---} & \cdots & \text{---} \\ \text{---} & w_k & \text{---} \end{bmatrix}$$

To classify an example: which row has the highest dot product with x ?

$$z = Wx + b$$

$$g(\hat{y}) = \operatorname{argmax}_k z_k$$

Converting to probabilities

We want probabilistic predictions: $p(y_k | x)$

i.e. $\sum_k p(y_k | x) = 1$, and $p(y_k | x) \geq 0$

Recall: use the softmax function:

$$\hat{y} = \text{softmax}(Wx + b)$$

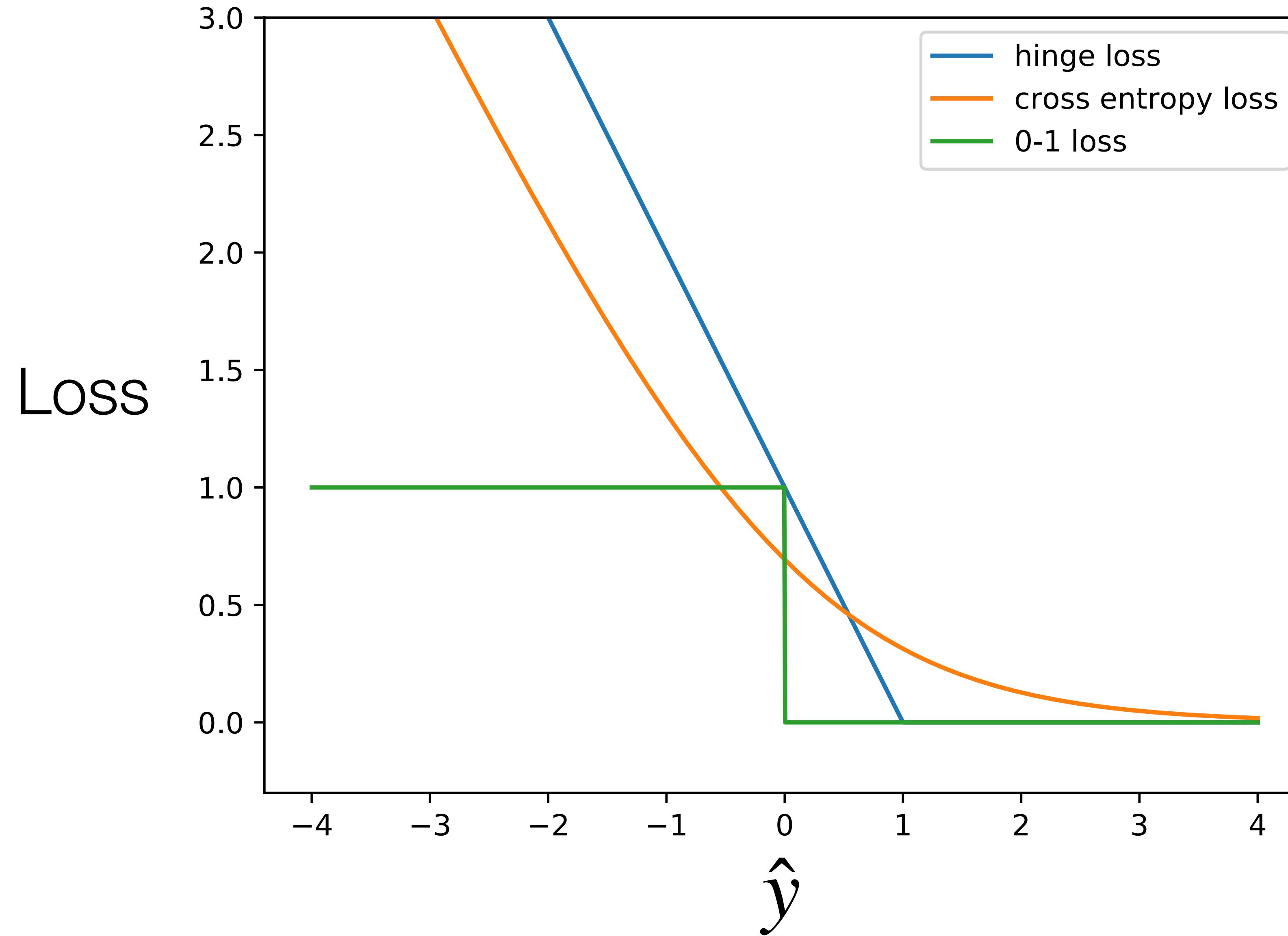
$$p(y_k | x) = \hat{y}_k$$

$$\text{softmax}(z)_k = \frac{\exp(z_k)}{\sum_i \exp(z_i)}$$

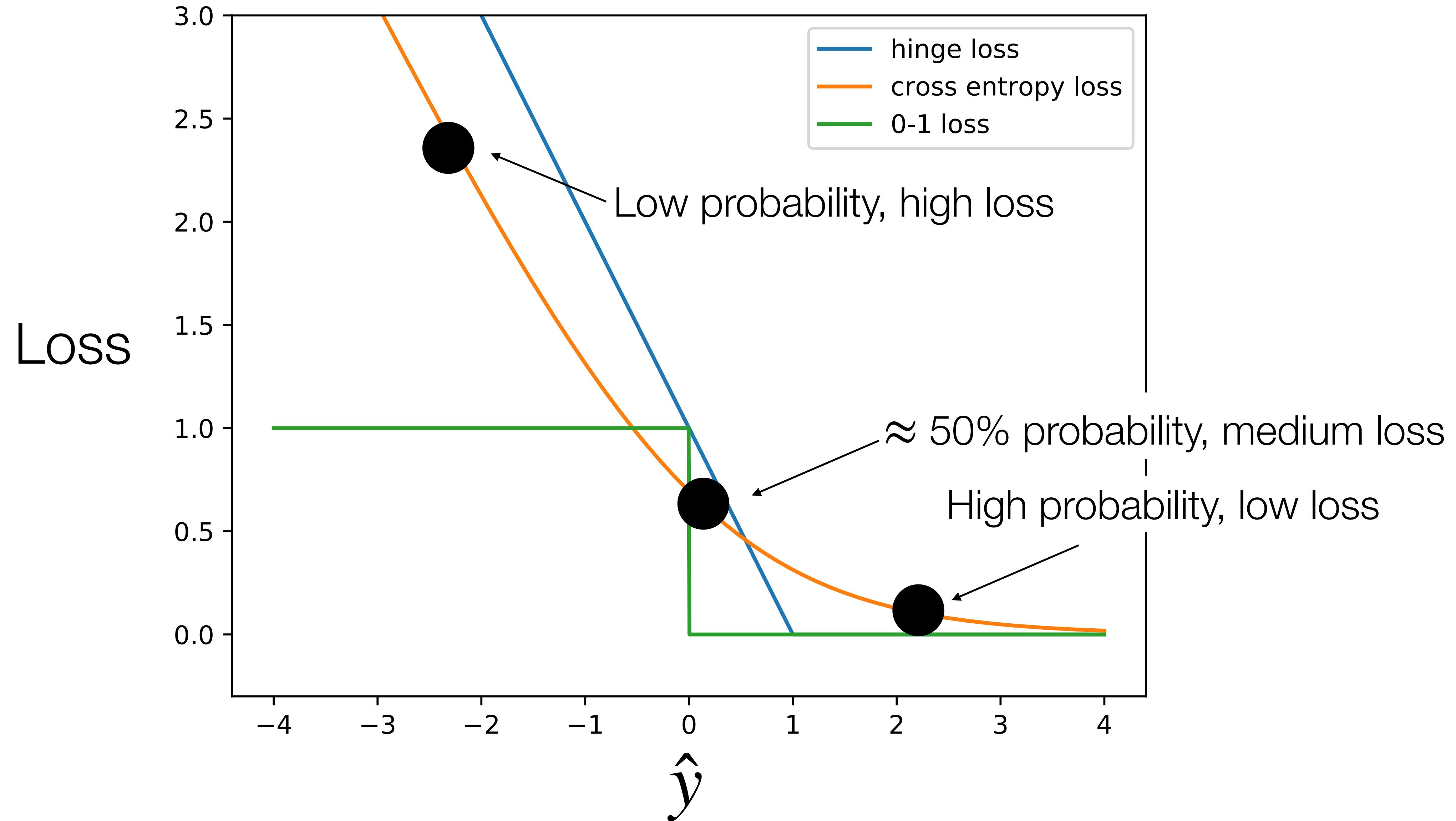
non-negative

sums-to-1

Loss functions



Loss functions

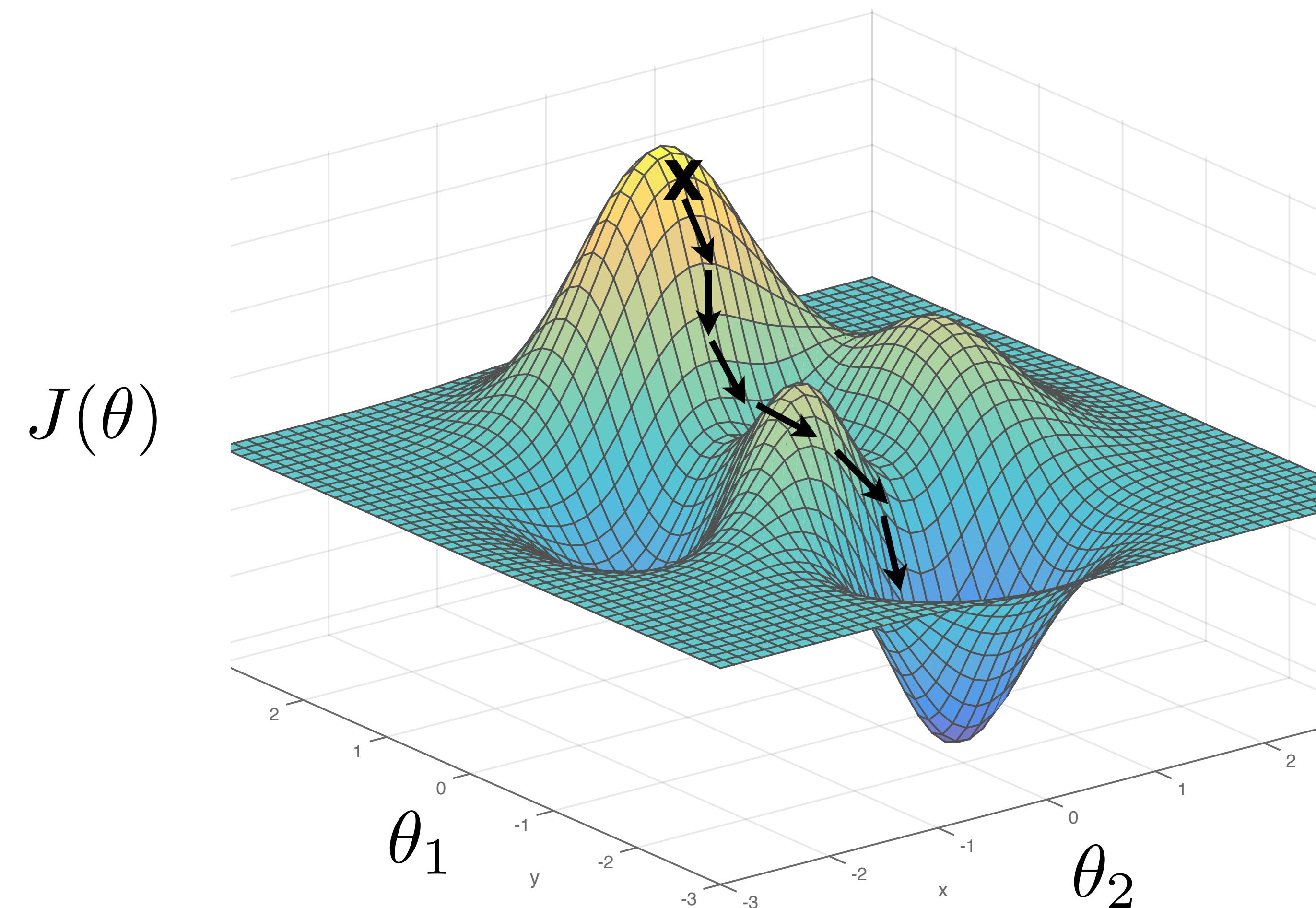


How do we learn a classifier?

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

$\overbrace{\hspace{10em}}$
 $J(\theta)$

Gradient descent



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

Gradient descent

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

$\overbrace{\hspace{10em}}$
 $J(\theta)$

One iteration of gradient descent:

$$\theta^{t+1} = \theta^t - \eta_t \nabla_{\theta} J(\theta) \quad \Big|_{\theta=\theta^t}$$

↓
learning rate

Gradient descent

$$\theta^{t+1} = \theta^t - \eta_t \nabla_{\theta} J(\theta) \Big|_{\theta=\theta^t}$$

What's this again?

1. Gradient of the loss w.r.t. the classifier's parameters
2. Direction of steepest descent
3. “Local” linear approximation to the function

Gradient descent

Take gradient at current θ

$$\theta^{t+1} = \theta^t - \eta_t \left. \nabla_{\theta} J(\theta) \right|_{\theta=\theta^t}$$

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\theta_1} \\ \frac{\partial J}{\theta_2} \\ \dots \\ \frac{\partial J}{\theta_d} \end{bmatrix}$$

Vector of partial derivatives for loss

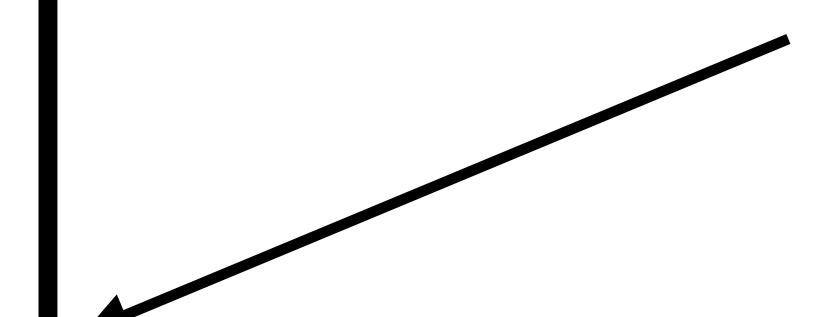
The diagram consists of two parts. On the left, there is a vector equation for the gradient. On the right, there is a text label. Two arrows point from the text label to the vector components. One arrow points from 'Vector of partial derivatives for loss' to the first term $\frac{\partial J}{\theta_1}$. Another arrow points from the same text label to the last term $\frac{\partial J}{\theta_d}$.

Estimating the gradient

Idea #1: finite differences $\frac{f(x + \epsilon) - f(x)}{\epsilon}$

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J}{\theta_1} \\ \frac{\partial J}{\theta_2} \\ \dots \\ \frac{\partial J}{\theta_d} \end{bmatrix} \approx \frac{J(\theta + \epsilon e_2) - J(\theta - \epsilon e_2)}{2\epsilon}$$

where e_2 is vector of zeros except for a 1 in 2nd component



Estimating the gradient

Idea #1: finite differences $\frac{f(x + \epsilon) - f(x)}{\epsilon}$

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J}{\theta_1} \\ \frac{\partial J}{\theta_2} \\ \dots \\ \frac{\partial J}{\theta_d} \end{bmatrix}$$

- Easy to compute but **slow**.
- Inaccurate in some cases, unless careful
- Useful for **checking** other methods

Estimating the gradient

Idea #2: compute it analytically using calculus.

Simple example: binary labels, squared loss, and regularization on θ

$$J(\theta) = \lambda \|\theta\|^2 + \frac{1}{N} \sum_{i=1}^N (y_i - \theta^\top x_i)^2$$

$$\nabla_\theta$$



$$\nabla_\theta J(\theta) = 2\lambda\theta - \frac{1}{N} \sum_{i=1}^N 2(y_i - \theta^\top x_i)x_i$$

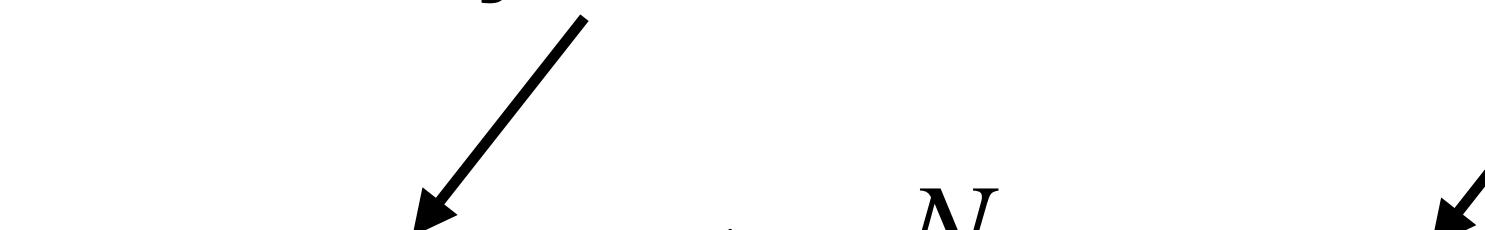
Analyzing the gradient descent update

Recall update rule: $\theta^{t+1} = \theta^t - \eta_t \nabla_{\theta} J(\theta)$ $|_{\theta=\theta^t}$

How does this change θ ?

$$-\nabla_{\theta} J(\theta) = -2\lambda\theta - \frac{1}{N} \sum_{i=1}^N \boxed{2(y_i - \theta^{\top} x_i)x_i}$$

“Decay” toward 0 Scalar α

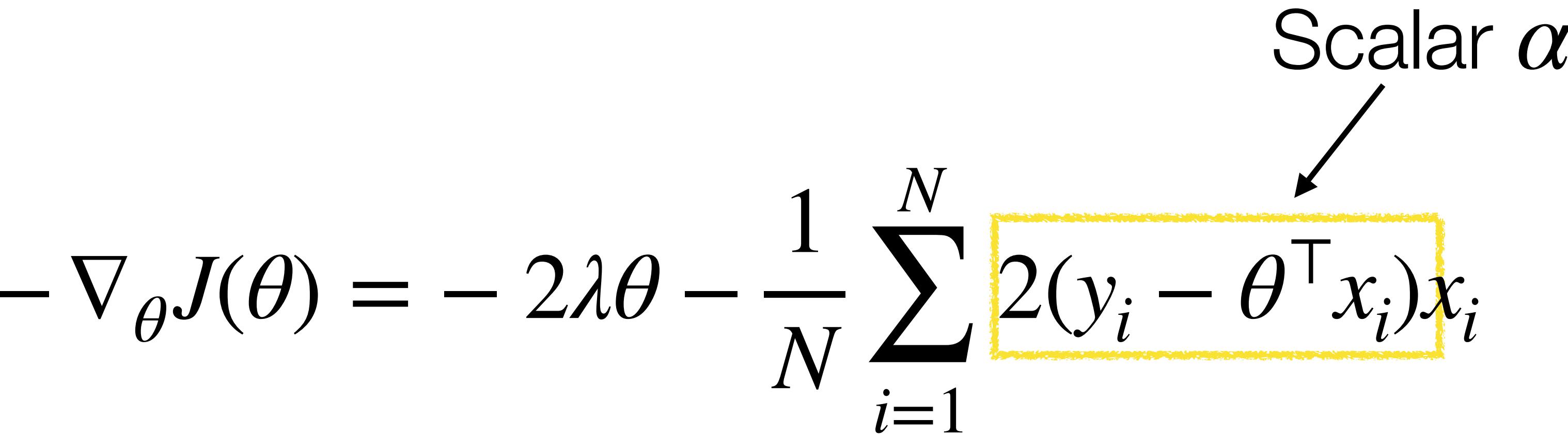


Analyzing the gradient descent update

What happens at each example?

$$-\nabla_{\theta} J(\theta) = -2\lambda\theta - \frac{1}{N} \sum_{i=1}^N \boxed{2(y_i - \theta^T x_i)x_i}$$

Scalar α



If $y > \theta^T x$ (too low): then $\theta_{t+1} = \theta + \alpha x$ for some α

Dot product before: $\theta^T x$

Dot product after: $(\theta + \alpha x)^T x = \theta^T x + \alpha x^T x$

Computational issues with gradient descent

Batch gradient descent

Loss function:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N L(x_i, y_i, \theta)$$

Its gradient is the sum of gradients for each example:

$$\nabla J(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla L(x_i, y_i, \theta)$$

Requires iterating over every training example each gradient step!

Can we speed this up?

Stochastic gradient descent

This is just an **average**:

$$\nabla J(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla L(x_i, y_i, \theta)$$

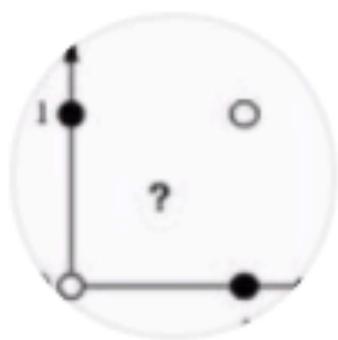
We know from statistics that we can estimate the average of a full “population” from a **sample**.

$$\nabla J(\theta) \approx \frac{1}{|B|} \sum_{i \in B} \nabla L(x_i, y_i, \theta)$$

where B is a **minibatch**: a random subset of examples.

This is called **stochastic gradient descent (SGD)**.

Stochastic gradient descent



ML Hipster @ML_Hipster · Jul 20, 2012

"Oh sure, going in that direction will totally minimize the objective function" —Sarcastic Gradient Descent.

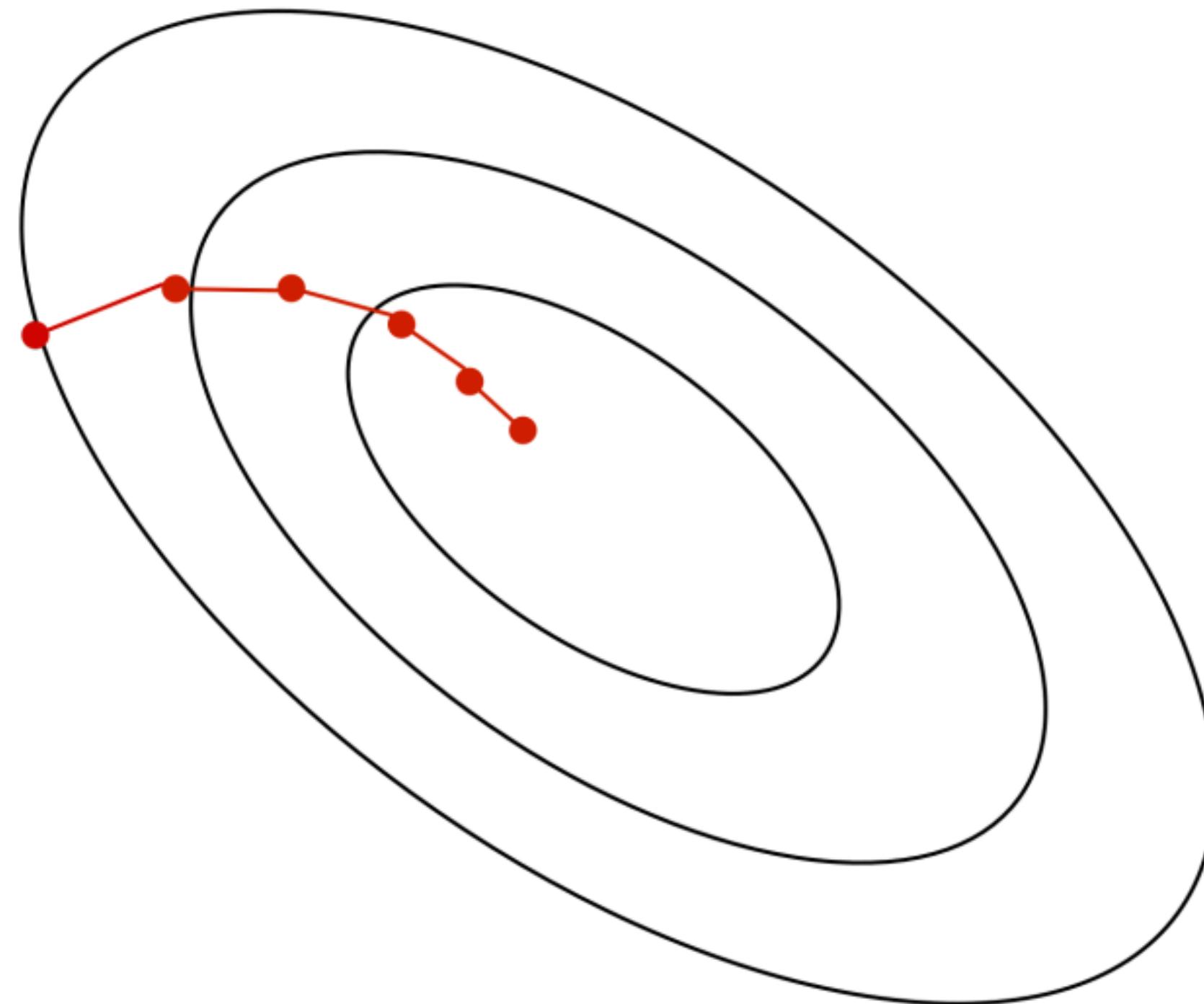
4

288

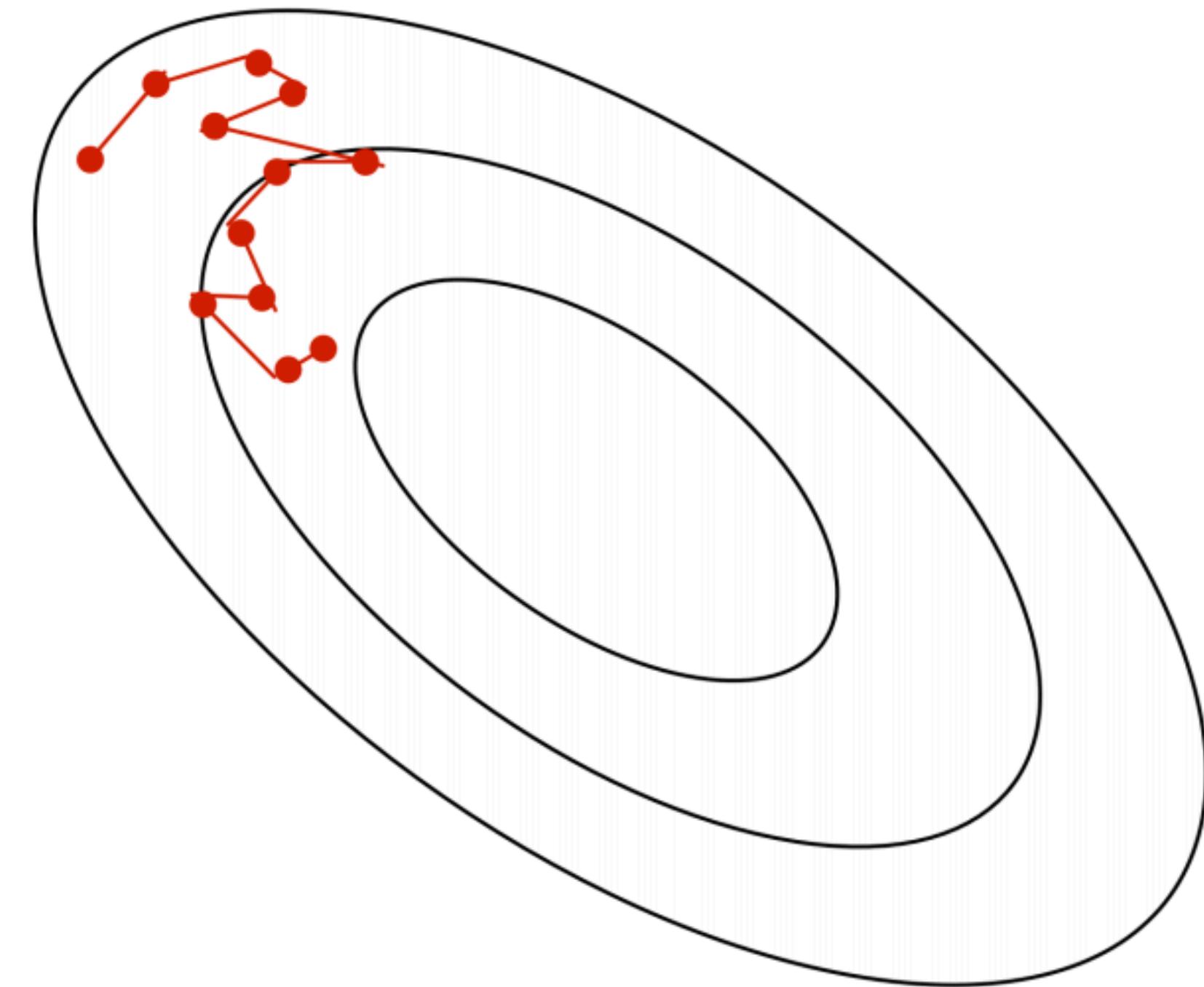
182



Stochastic gradient descent



Batch gradient descent

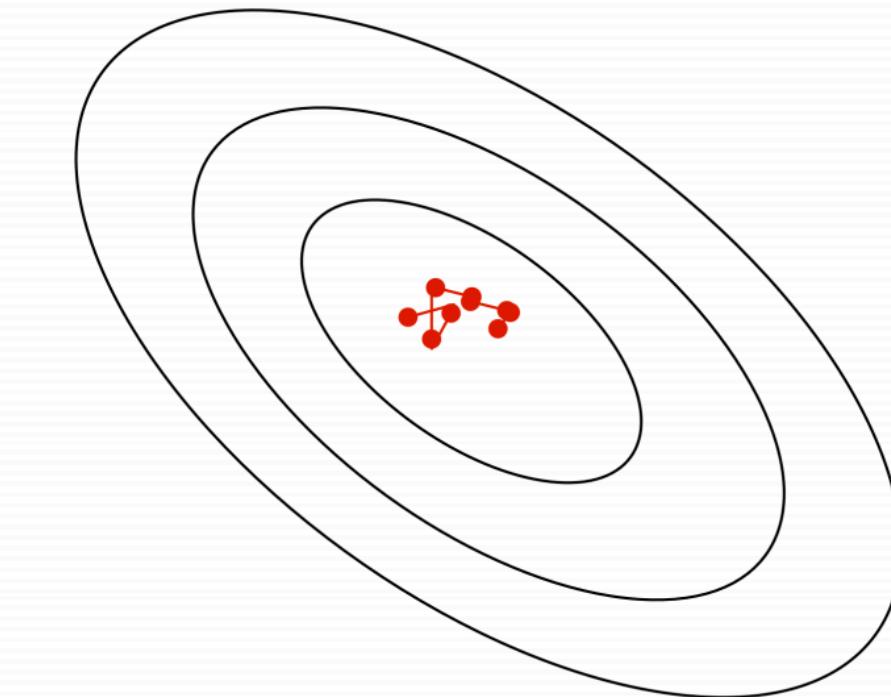


Stochastic gradient descent

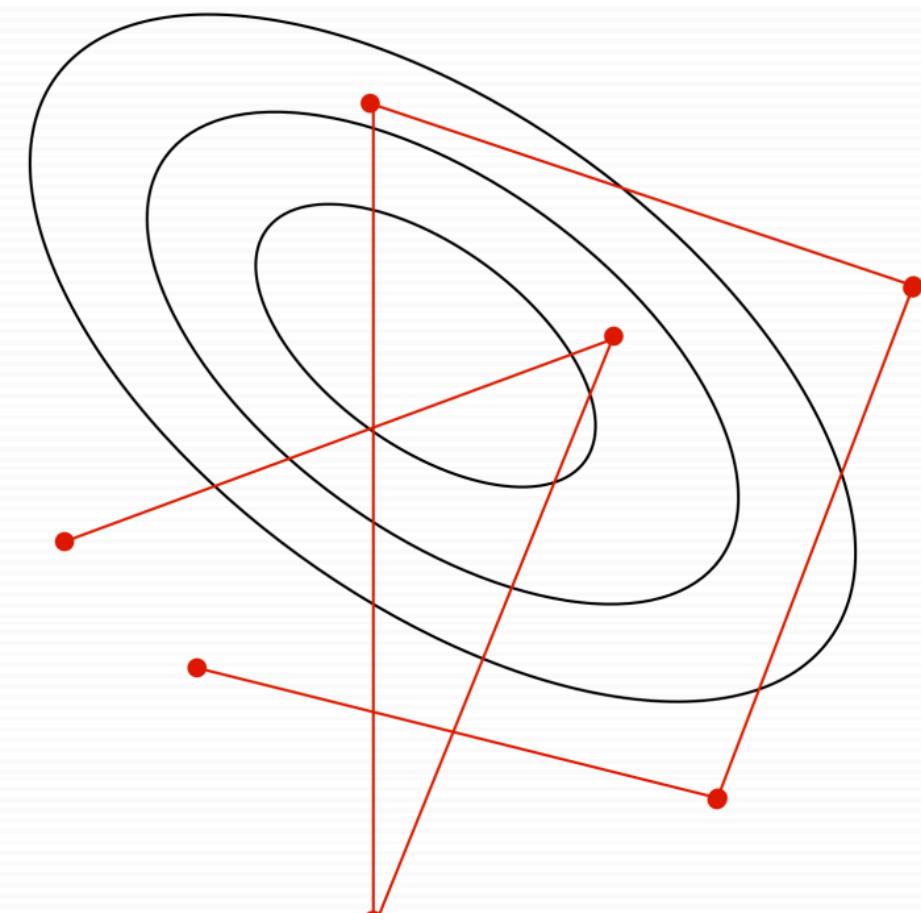
Learning rate

- Sensitive to the learning rate:

$$\theta^{t+1} = \theta^t - \eta_t \nabla_{\theta} J(\theta) \Big|_{\theta=\theta^t}$$



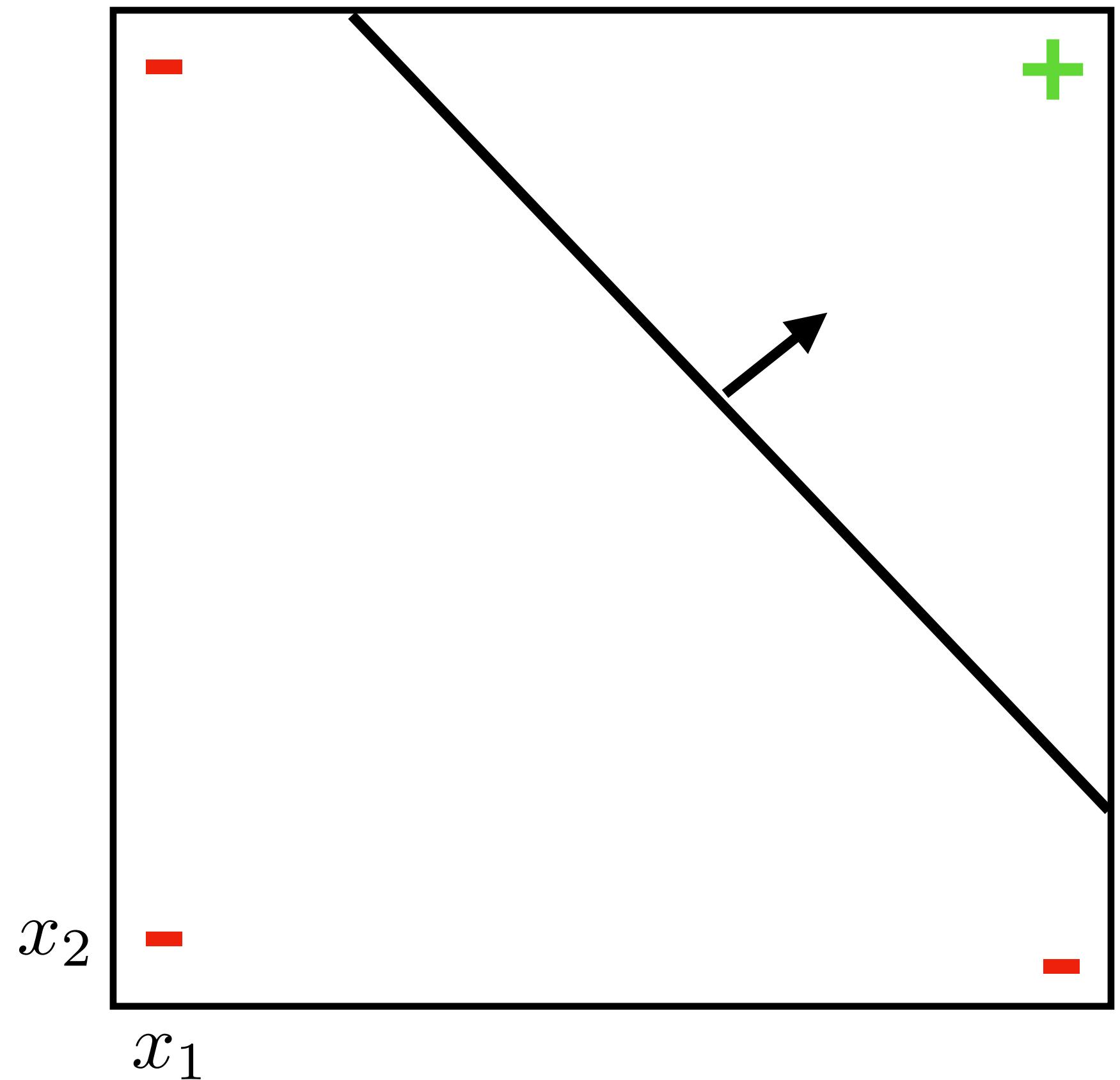
Small learning rate



Large learning rate

- Often start with high learning rate, and decrease over time
- For example: start with $\eta_t = 0.1$, then decrease to $\eta_t = 0.01$ when training loss plateaus

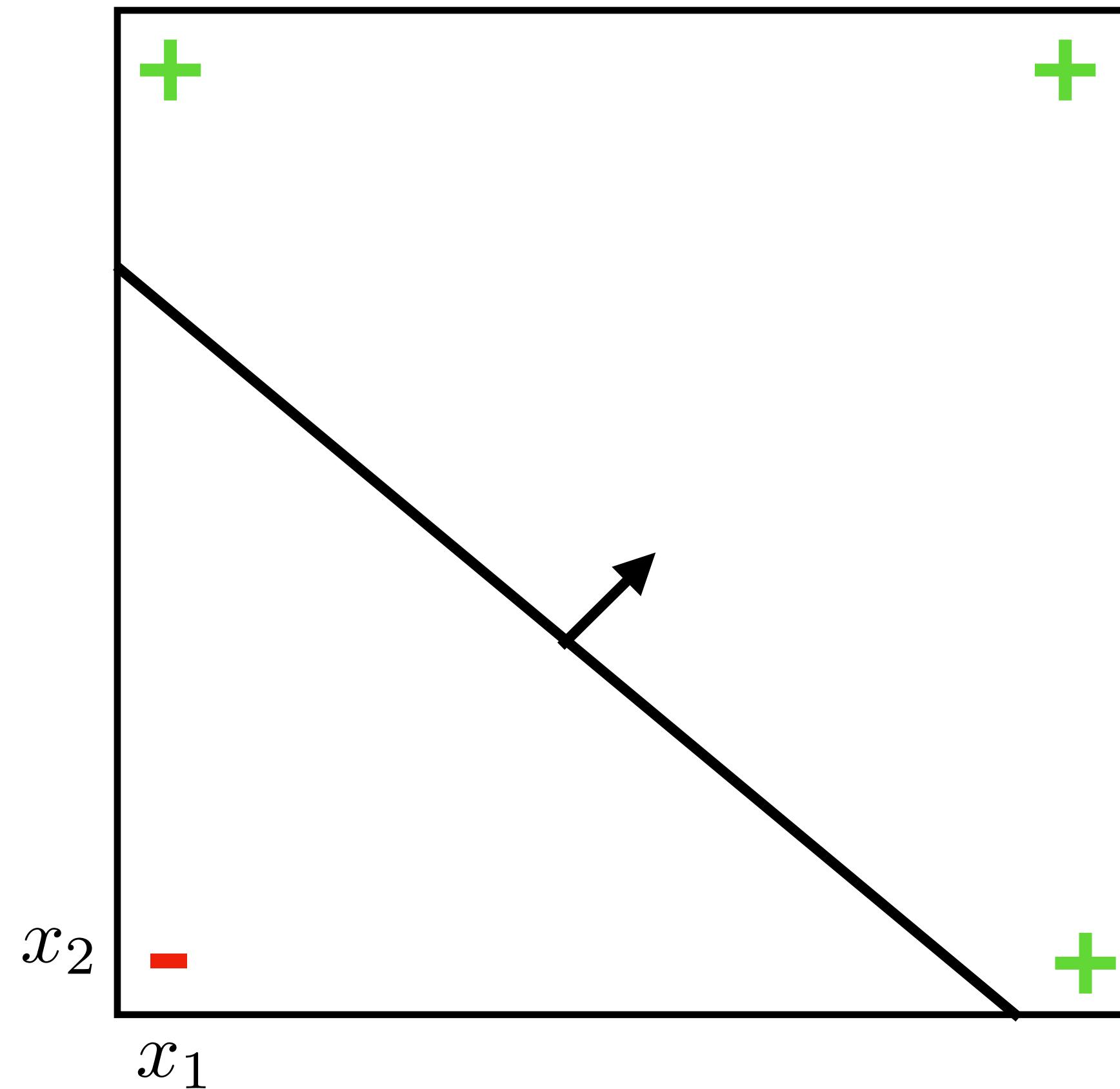
Limitations to linear classifiers



		x_2
		0
x_1	0	1
0	0	0
1	0	1

AND

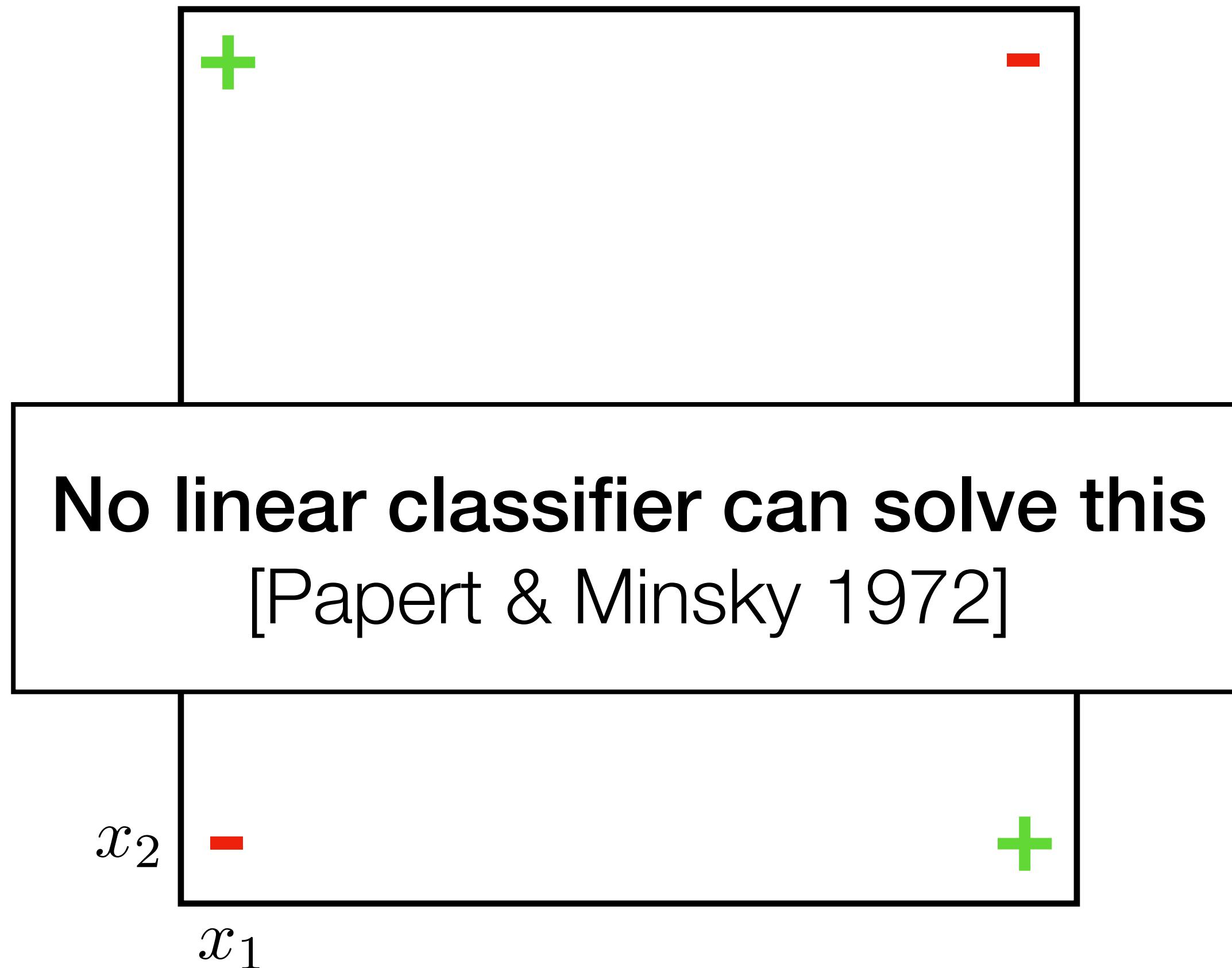
Limitations to linear classifiers



		x_2	
		0	1
x_1	0	0	1
	1	1	1

OR

Limitations to linear classifiers



	x_2	
x_1	0	1
	0	0
	1	0
		1

XOR

Next lecture: neural networks