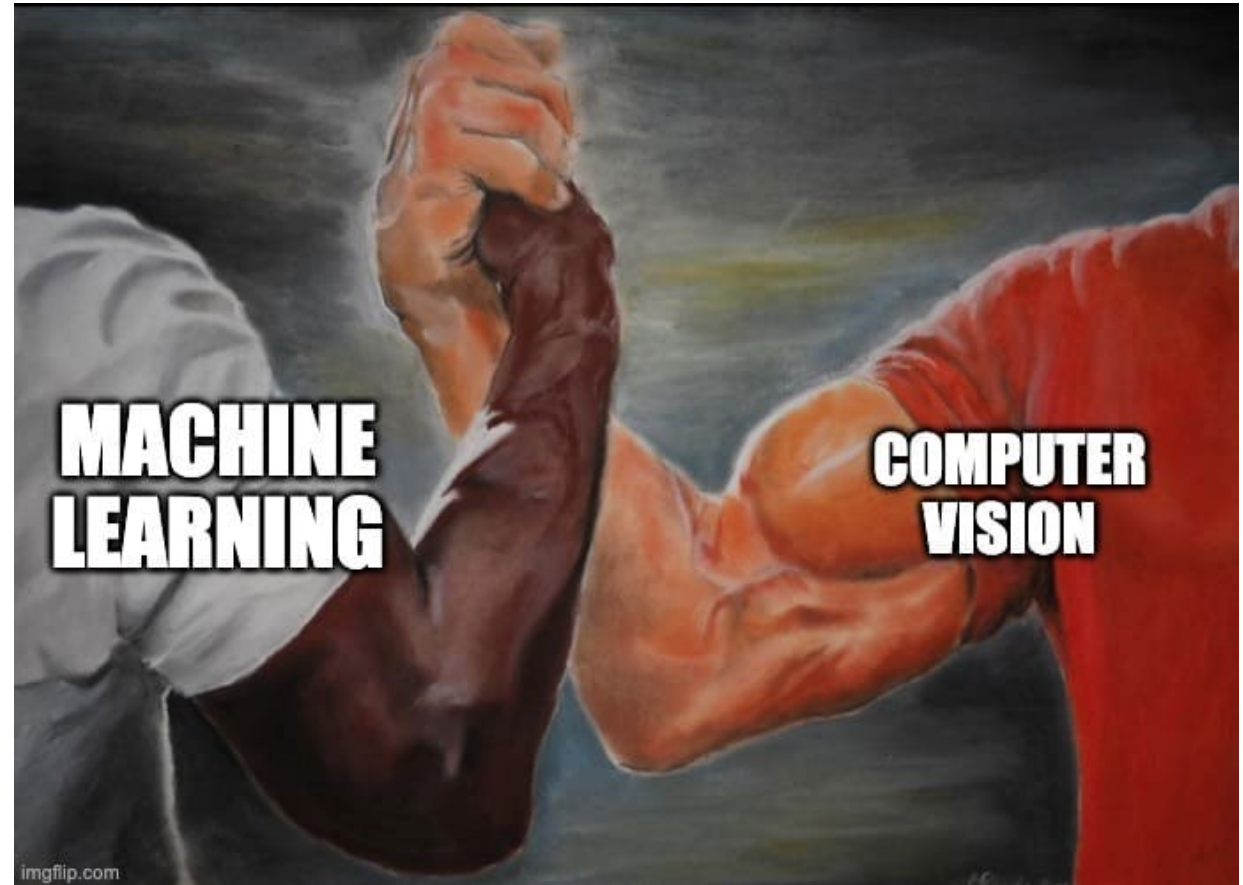


Lecture 7

Machine Learning for Computer Vision



Motivation: Image Classification



What is this?

{dog, cat, airplane, bus, laptop,
chair ...}

What animal is this ?

{dog, cat, lion, tiger, duck, cow,
giraffe, ...}

What type of cat is this?

{Cheshire, Siamese, Persian,
Shorthair, Bombay, ...}

Motivation: Image Classification



Central question: What “category” is this?

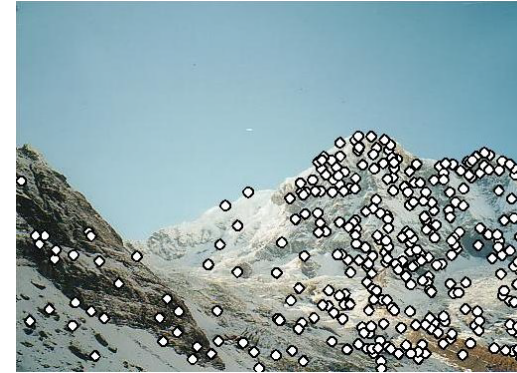
How can a computer vision system make a decision like this?

Ideas

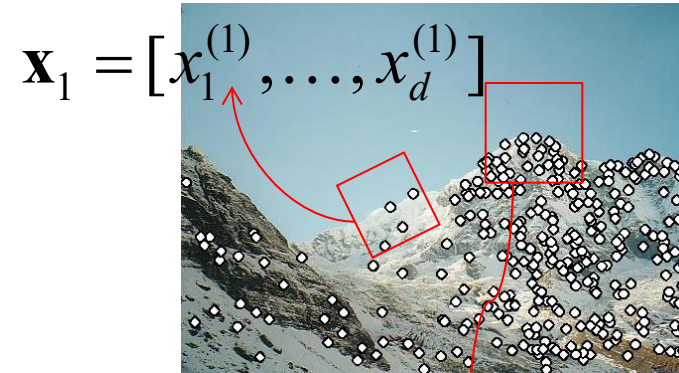
- Based on colors, textures, shapes, edges, ...
- Based on features !!!

Features: main components

1) Detection: Identify the interest points

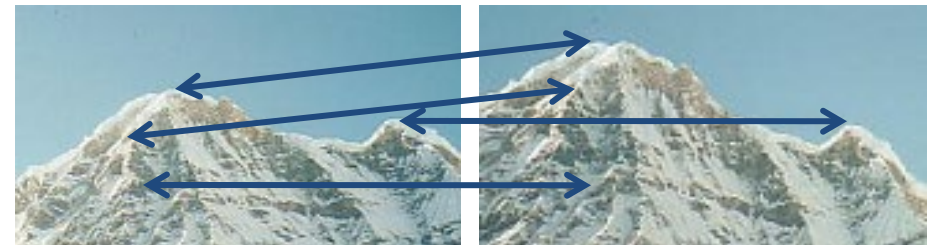


2) **Description: Extract vector feature descriptor surrounding each interest point.**



$$\mathbf{x}_2 = [x_1^{(2)}, \dots, x_d^{(2)}]$$

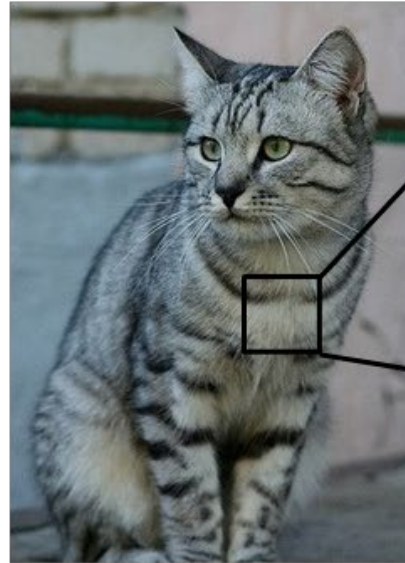
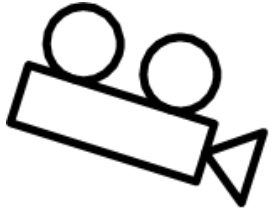
3) Matching: Determine correspondence between descriptors in two views



Recap & Motivation

- Image features are “interesting”, “unique” regions in an image
 - Intuitively these are “important”
- So far – we have seen how to detect and describe (*a.k.a. “represent”*) certain types of feature
 - Harris Corners, Blobs, ...
- We had a definition for what a “feature” is
 - Can we learn that from data?

Challenges: Viewpoint Variation

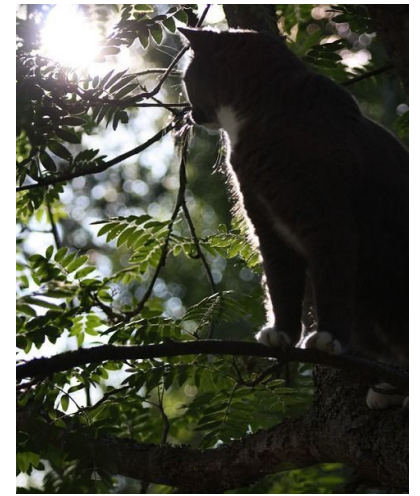


```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 98 105 128 105 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]  
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]  
[128 137 144 140 109 95 86 78 62 65 63 63 68 73 86 101]  
[125 133 148 137 110 121 117 94 65 79 80 65 54 64 72 98]  
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]  
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]  
[ 89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 80]  
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]  
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]  
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]  
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]  
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]  
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]  
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]  
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]  
[128 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]  
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]  
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]  
[122 164 148 103 71 56 78 83 103 119 139 102 61 69 84]]
```



All pixels change when
the camera moves!

Challenges: Illumination



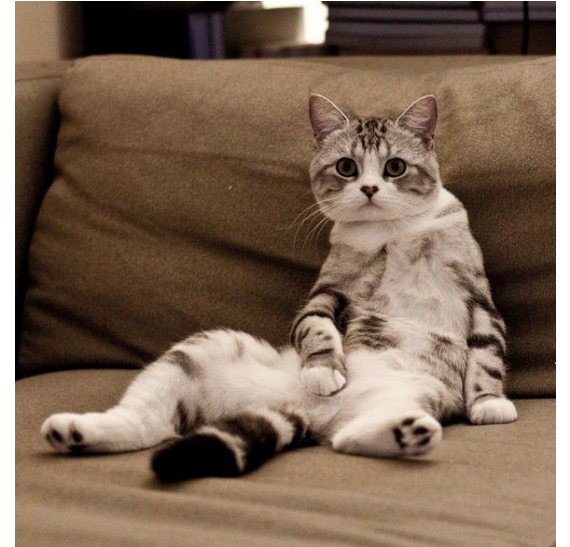
Challenges: Background Clutter



Challenges: Occlusion



Challenges: Pose and Deformation (Cat Yoga)

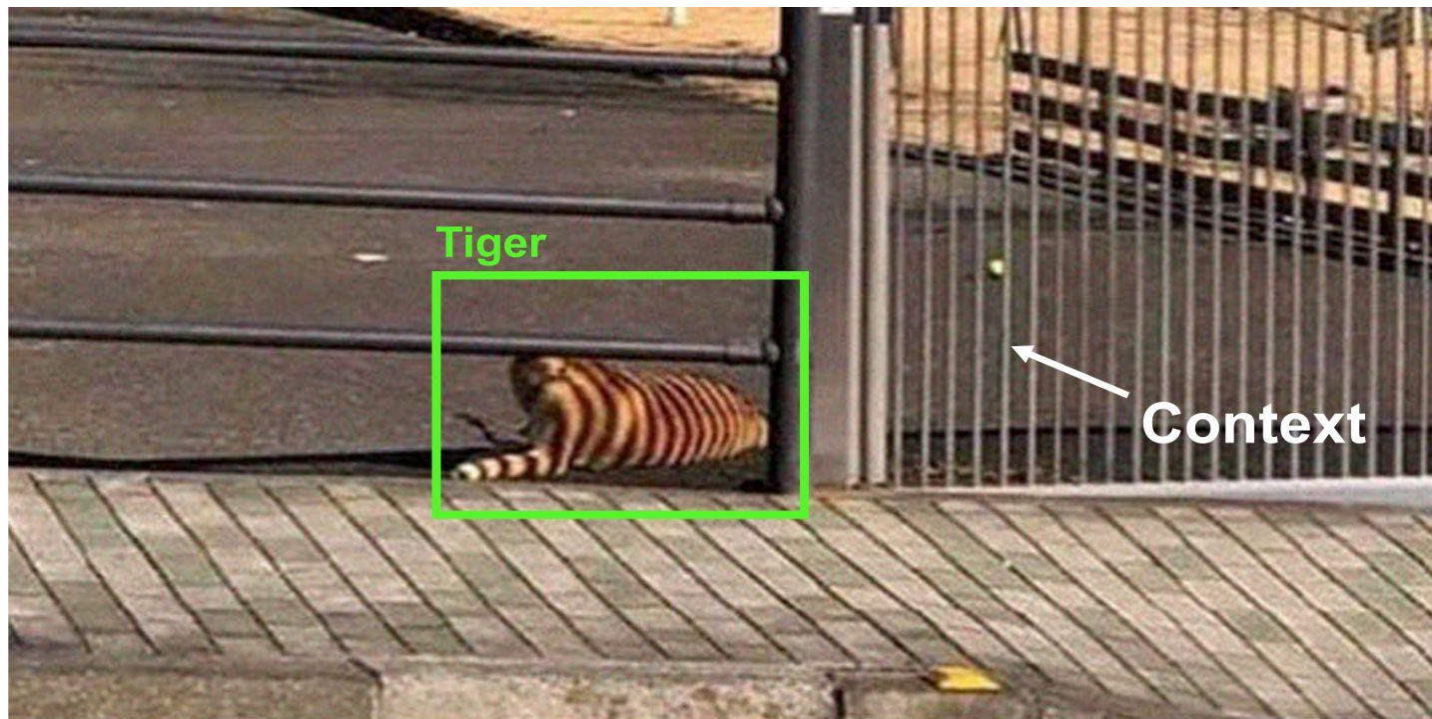


Challenges: Inter-Class Variation



Slide Credit: Fei-Fei Li

Challenges: Illusions



Data !!!



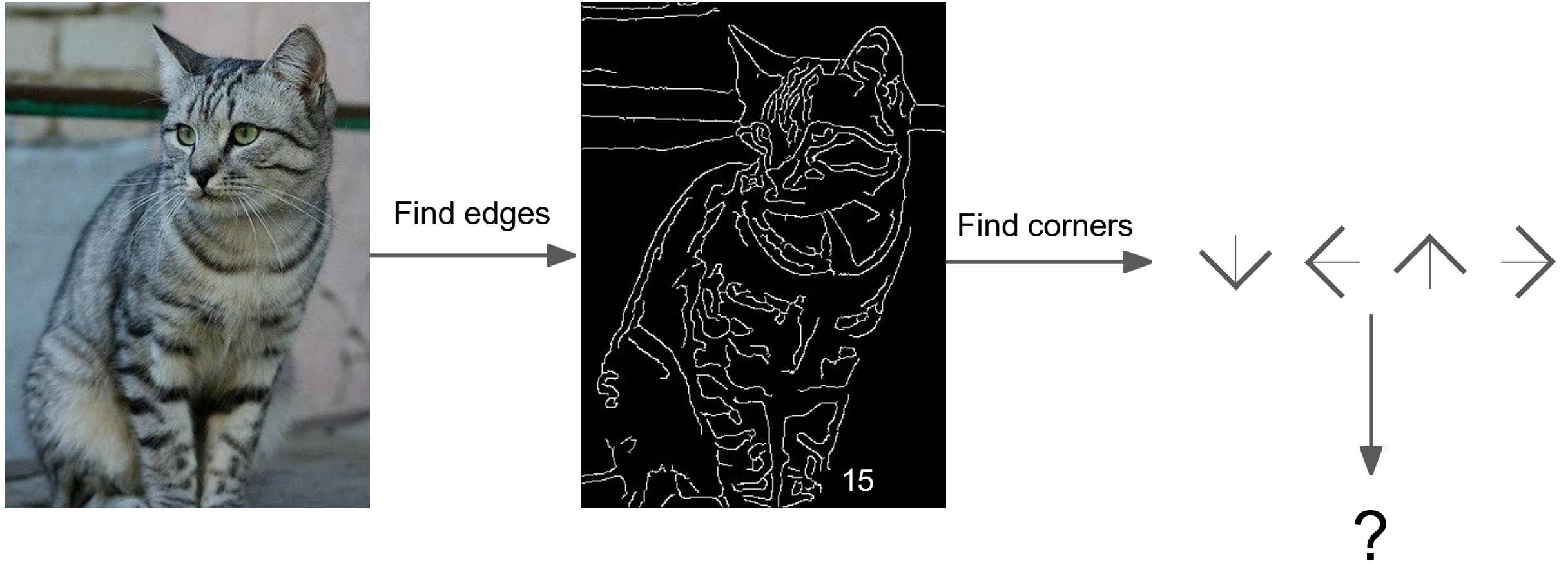
An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

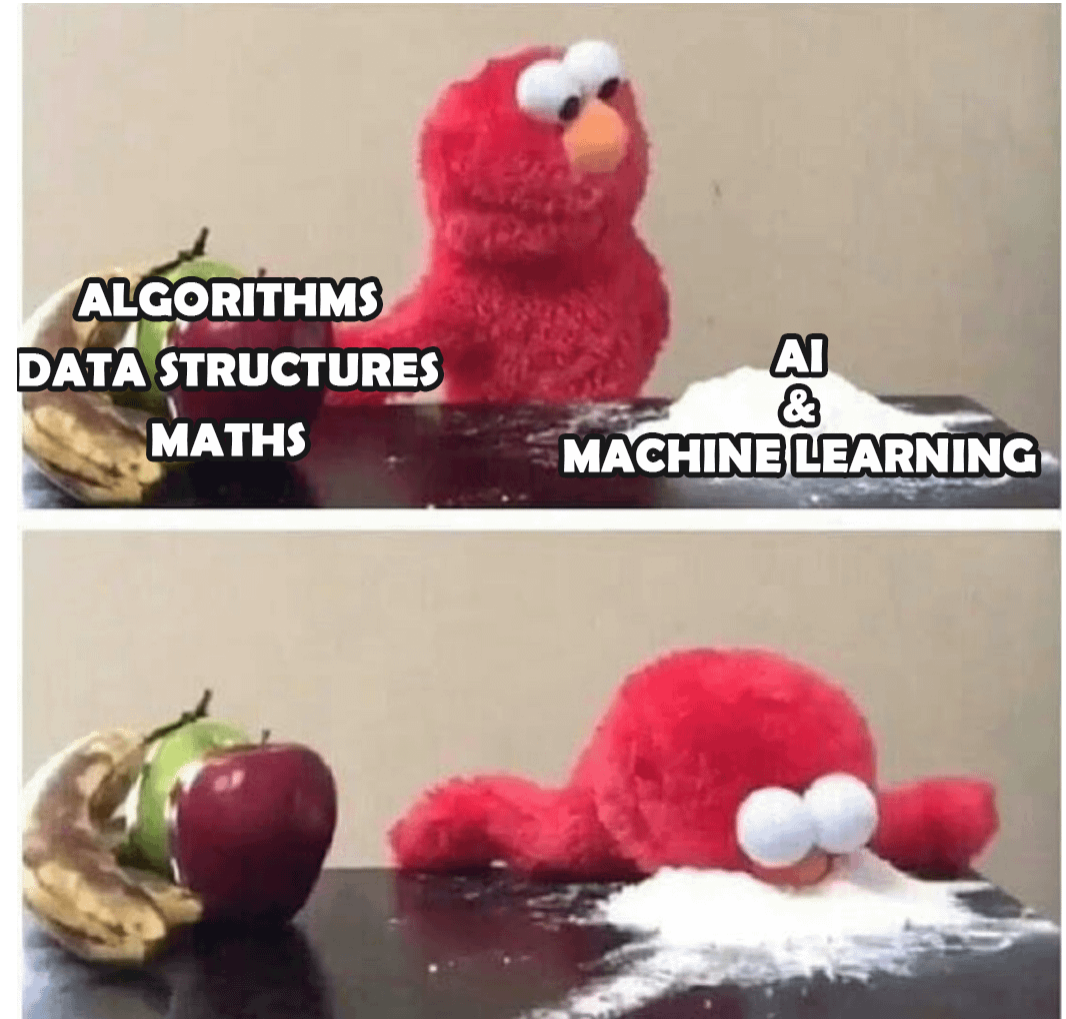
Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.

Can we use features to make the decision?

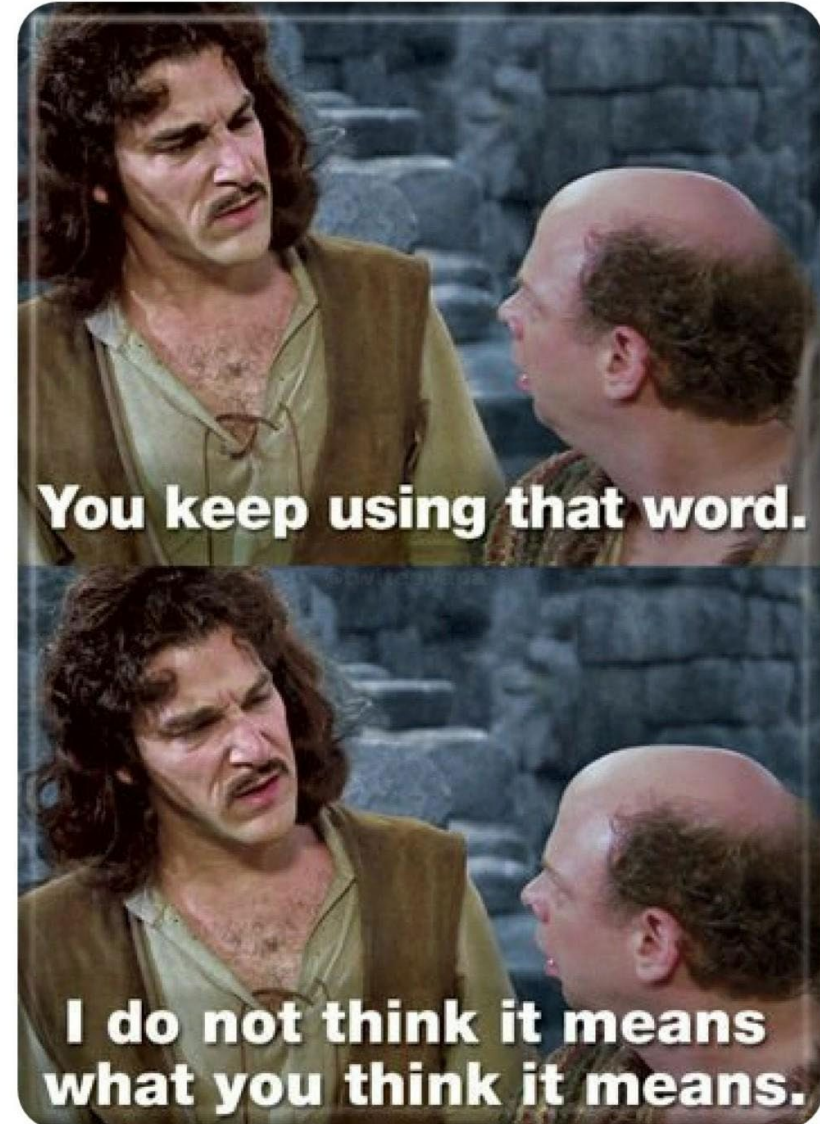


Get Ready for
The Good Stuff

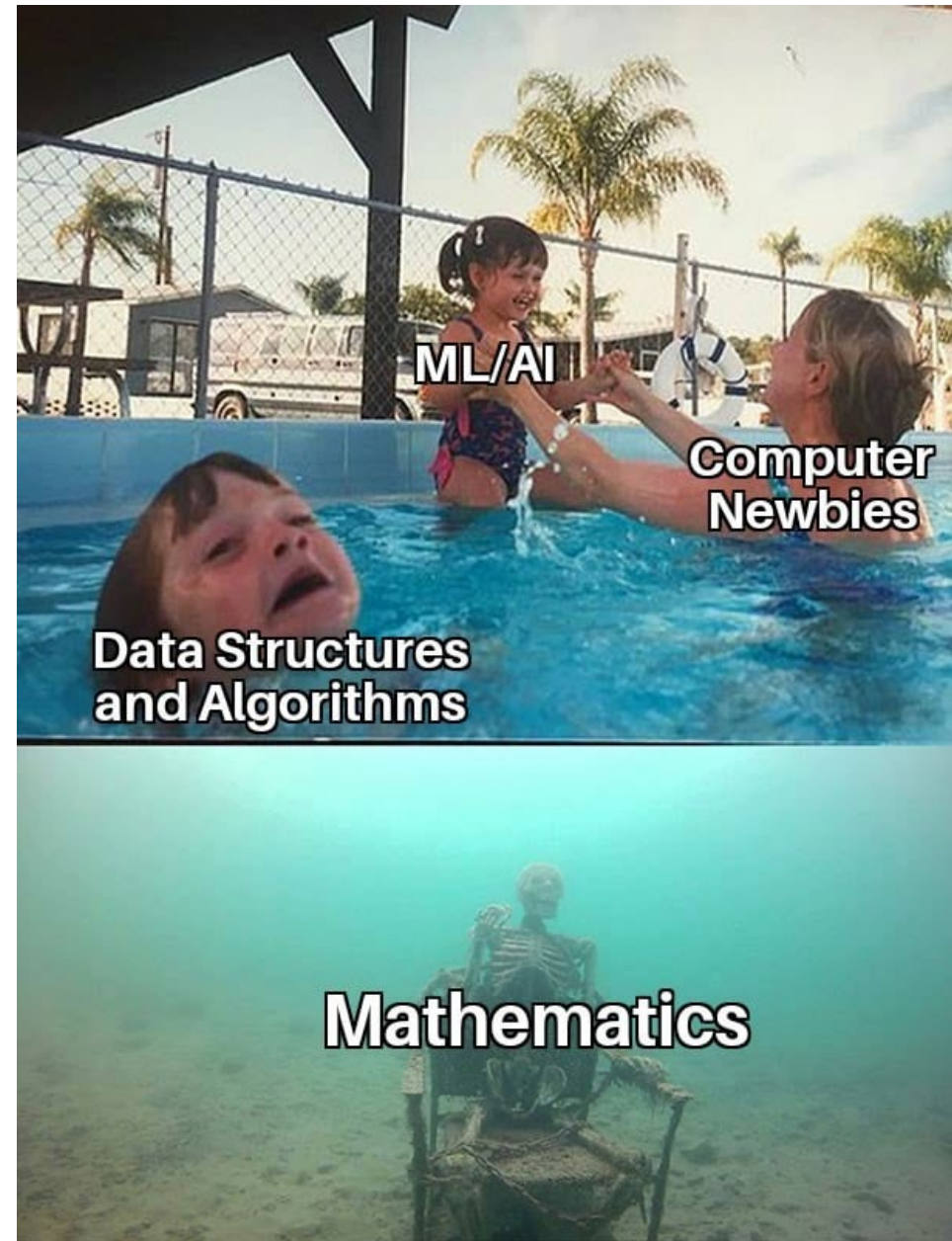


When someone uses 'Machine learning', 'AI' and 'deep learning' interchangeably in a discussion

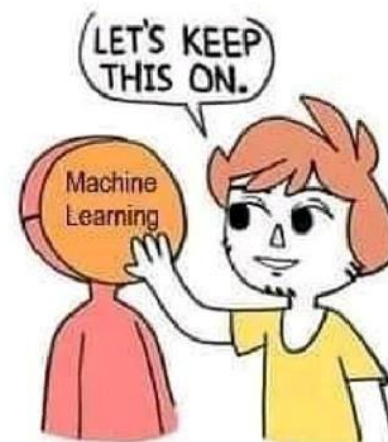
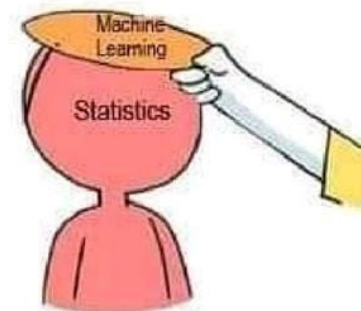
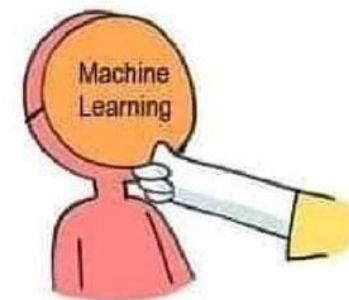
Wall Street / Silicon Valley
can you please stop



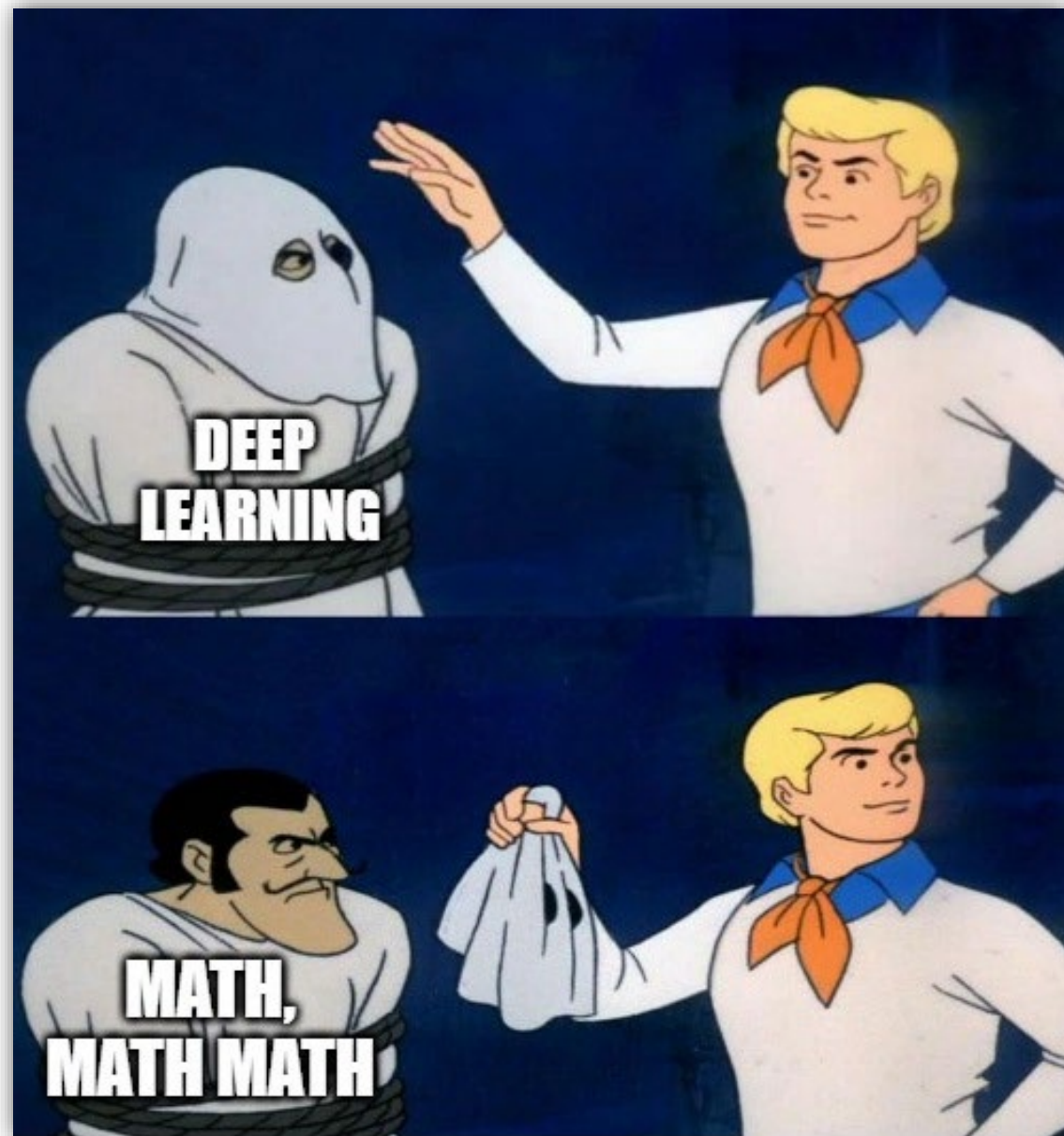
Know
your ancestors



The Open Secret

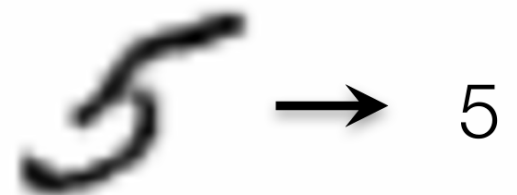
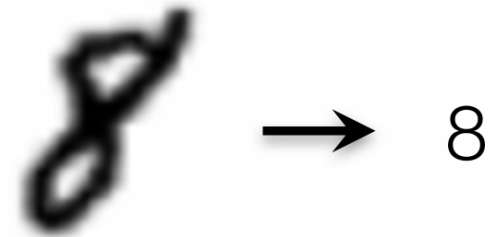


The Open Secret




“Learn” ?

- Let's look at a “programming” task
- The task:
Write a program that outputs the number in a 28x28 grayscale image




“Learn” ?

- **Approach 1:** try to write a program by hand
 - How would you do it ?



→ 8



→ 5

“Learn” ?

- **Approach 1:** try to write a program by hand
 - How would you do it ?
- **Approach 2:** (the machine learning approach)
 - Collect a large “dataset” of digit images
 - “Label” them with the corresponding numbers (0, 1, ..., 9)
 - Let the system “write its own program” to map from images to numbers

Machine Learning

1. Collect a dataset of images and labels
2. Use Machine Learning algorithms to train a classifier
3. Evaluate the classifier on new images

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

Example training set

airplane



automobile



bird



cat



deer



An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.

Learning

Data

```
def train(images, labels):  
    # Machine learning!  
    return model
```

Inference

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

→ Output

Nearest Neighbor Classifier

```
def train(images, labels):  
    # Machine learning!  
    return model
```



Memorize all
data and labels

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```



Predict the label
of the most similar
training image

Nearest Neighbor Classifier

deer



bird



plane



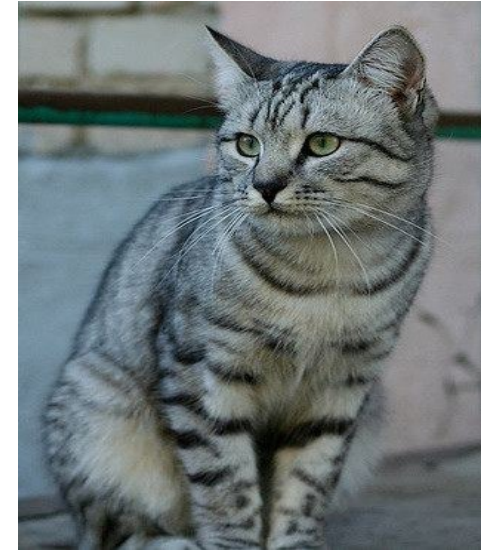
cat



car



Training data with labels



query data



Distance Metric



,



$\rightarrow \mathbb{R}$

Distance Metric to compare images

L1 distance:

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

test image

56	32	10	18
90	23	128	133
24	26	178	200
2	0	255	220

training image

10	20	24	17
8	10	89	100
12	16	178	170
4	32	233	112

-

=

pixel-wise absolute value differences

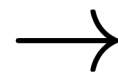
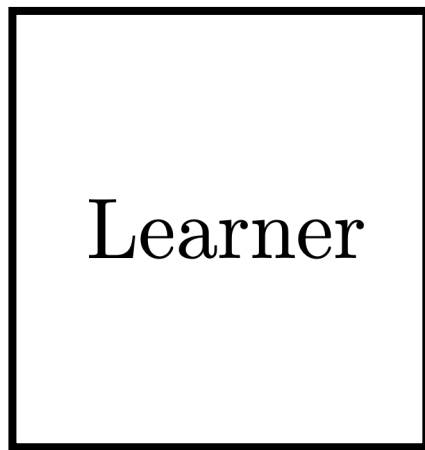
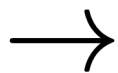
46	12	14	1
82	13	39	33
12	10	0	30
2	32	22	108

→ 456

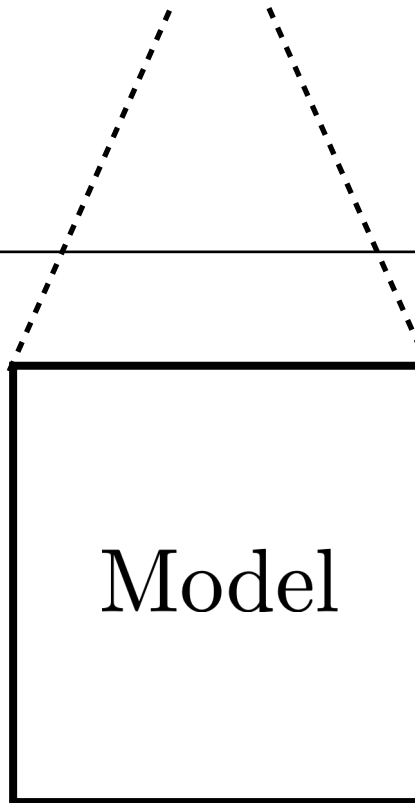
The goal of learning is to extract lessons from past experience in order to solve future problems.

Learning

Data

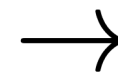
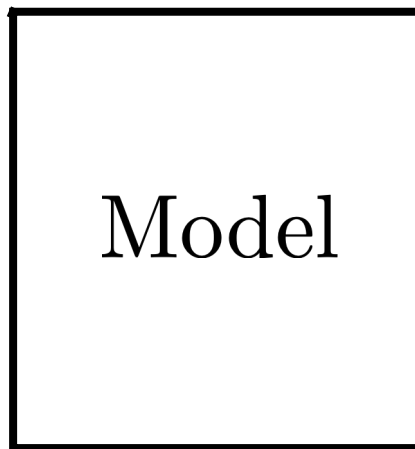
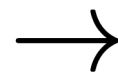


Model



Inference

Input



Output

Learning

Data

```
def train(images, labels):  
    # Machine learning!  
    return model
```

Inference

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

→ Output

The goal of learning is to extract lessons from past experience in order to solve future problems.

Let's LEARN. What does ☆ do?

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

Goal: answer future queries involving ☆

Approach: figure out what ☆ is doing by observing its behavior on examples

Past experience

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

Future query

$$3 \star 5 = ?$$

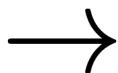
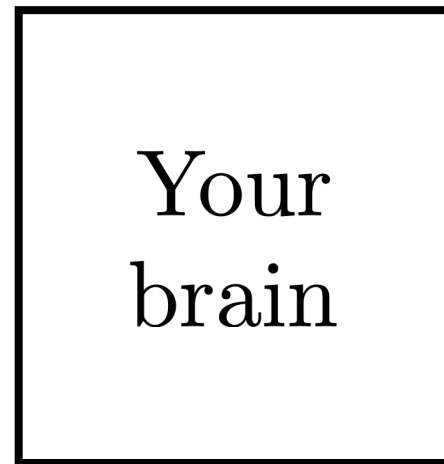
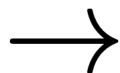
Training

$$2 \star 3 = 36$$

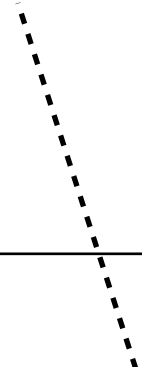
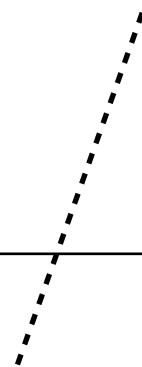
$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

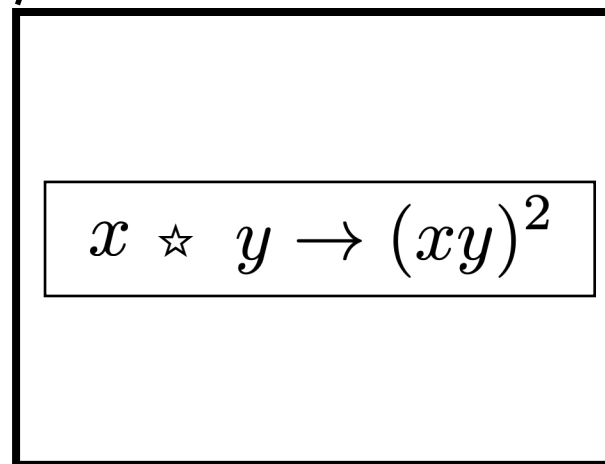
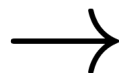


$$x \star y \rightarrow (xy)^2$$



Testing

$$3 \star 5$$



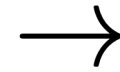
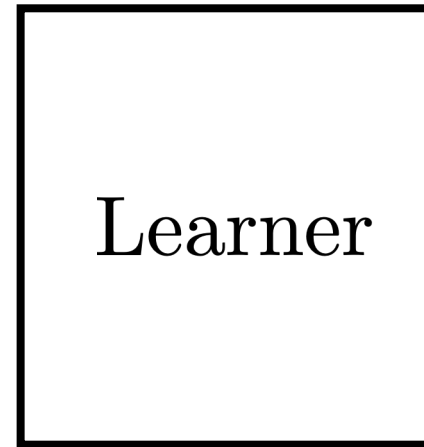
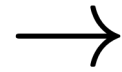
$$\rightarrow 225$$

Learning from examples

(aka **supervised learning**)

Training data

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



f

The goal of learning is
to extract lessons from past experience
in order to solve future problems.

Learning from examples

(aka **supervised learning**)

Training data

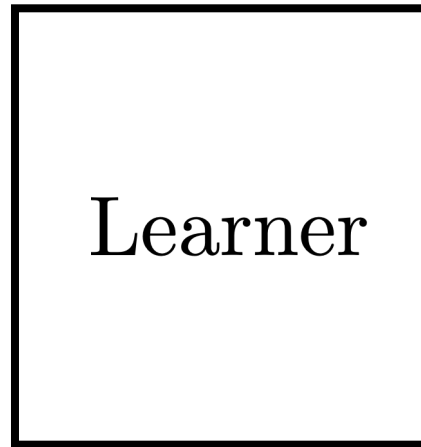
$\{x^{(1)}, y^{(1)}\}$

$\{x^{(2)}, y^{(2)}\}$

$\{x^{(3)}, y^{(3)}\}$

...

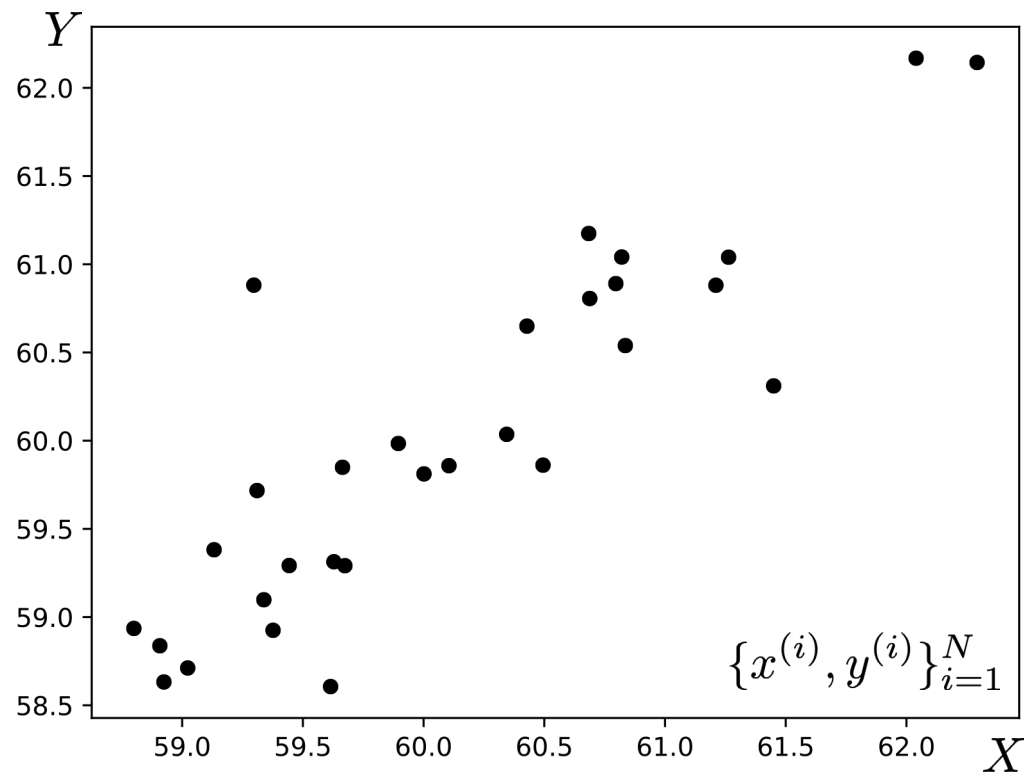
\rightarrow



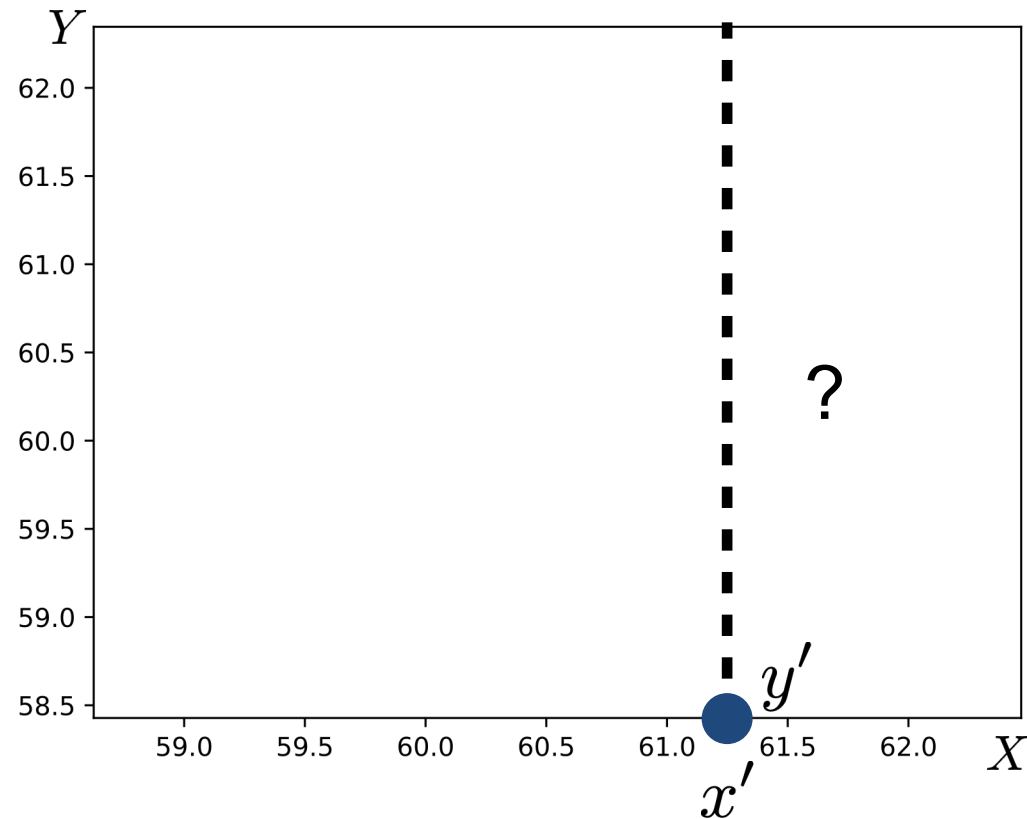
\rightarrow

$f : X \rightarrow Y$

Training data



Test query



Real-World Application:

A Model for Predicting Electricity Use

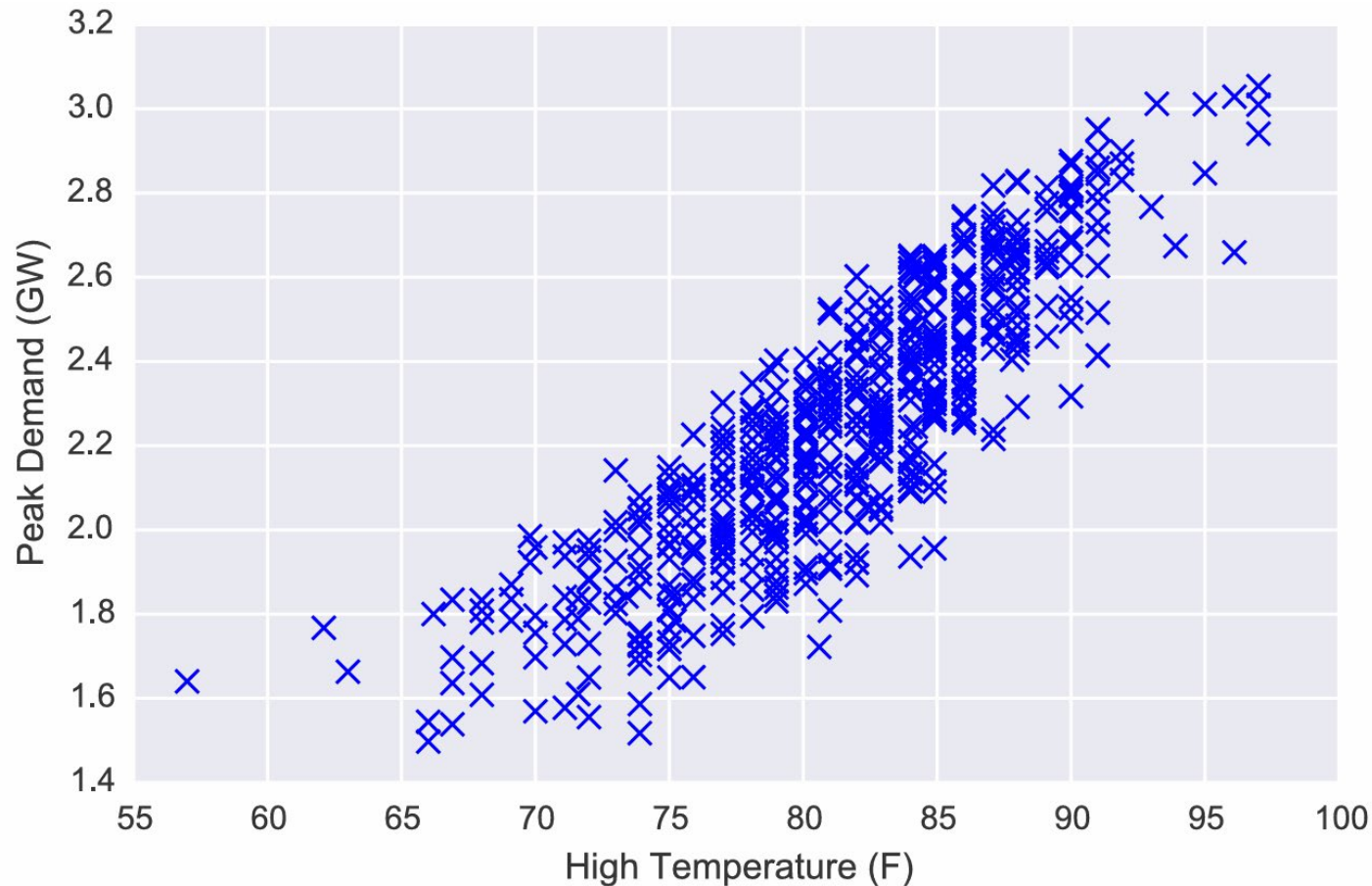
- What will the peak power consumption be in <your-favorite-city> tomorrow?
- Difficult to answer this question without data
 - Difficult to build an “a priori” model from first principles ...
- Relatively easy to record consumption history (the utility company has this data)
- Relatively easy to record features that may affect consumption:
 - temperature

Date	High Temperature (F)	Peak Demand (GW)
2011-06-01	84.0	2.651
2011-06-02	73.0	2.081
2011-06-03	75.2	1.844
2011-06-04	84.9	1.959
...

Real-World Application:

A Model for Predicting Electricity Use

- What will the peak power consumption be in <your-favorite-city> tomorrow?

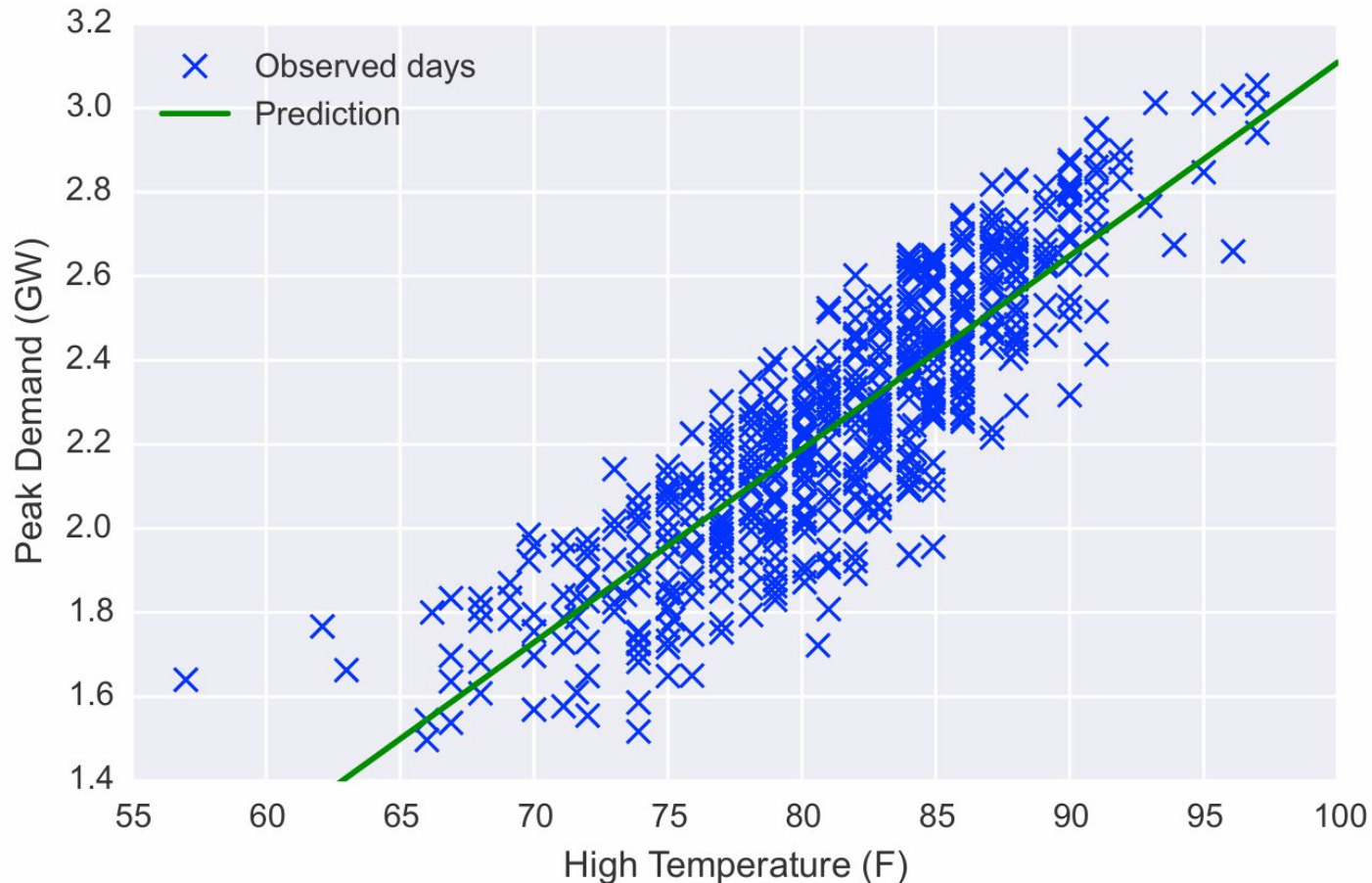


← DATA

Real-World Application:

A Model for Predicting Electricity Use

- What will the peak power consumption be in <your-favorite-city> tomorrow?

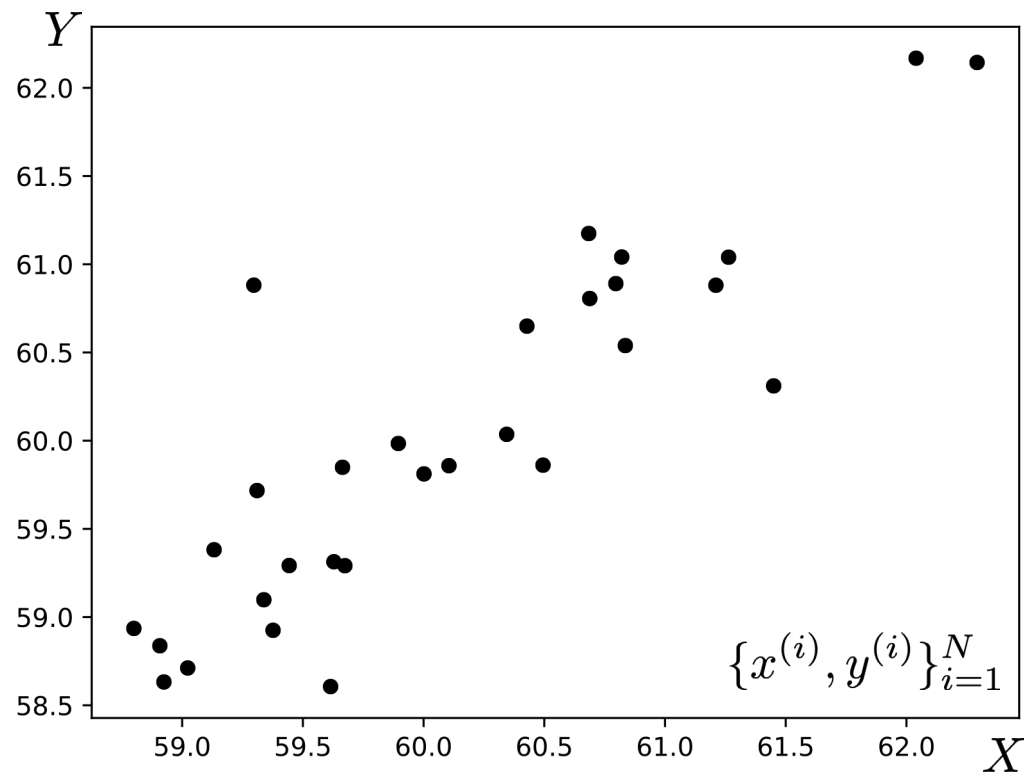


← PREDICTION

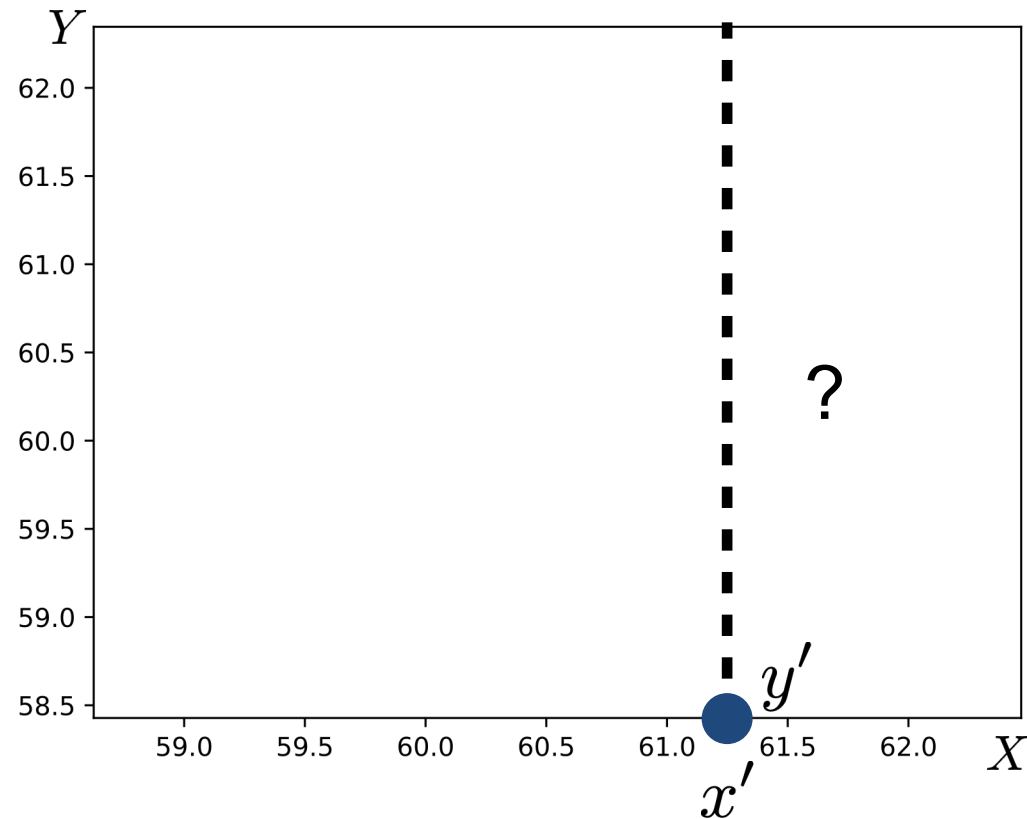
The essence of machine learning:

- A pattern exists
- We cannot pin down the pattern as an equation
- We need to approximate the pattern as a function of the input
 - Using data!
 - For vision: the pattern is in terms of pixel intensities and features.

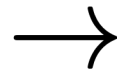
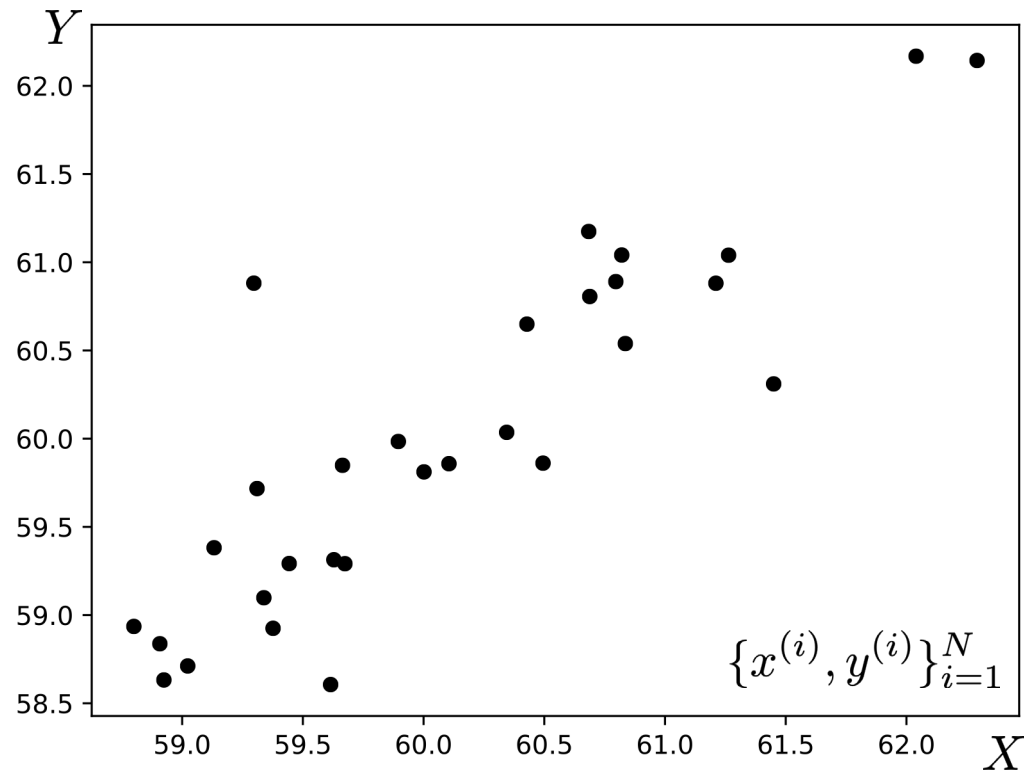
Training data



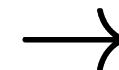
Test query



Training data



Learner

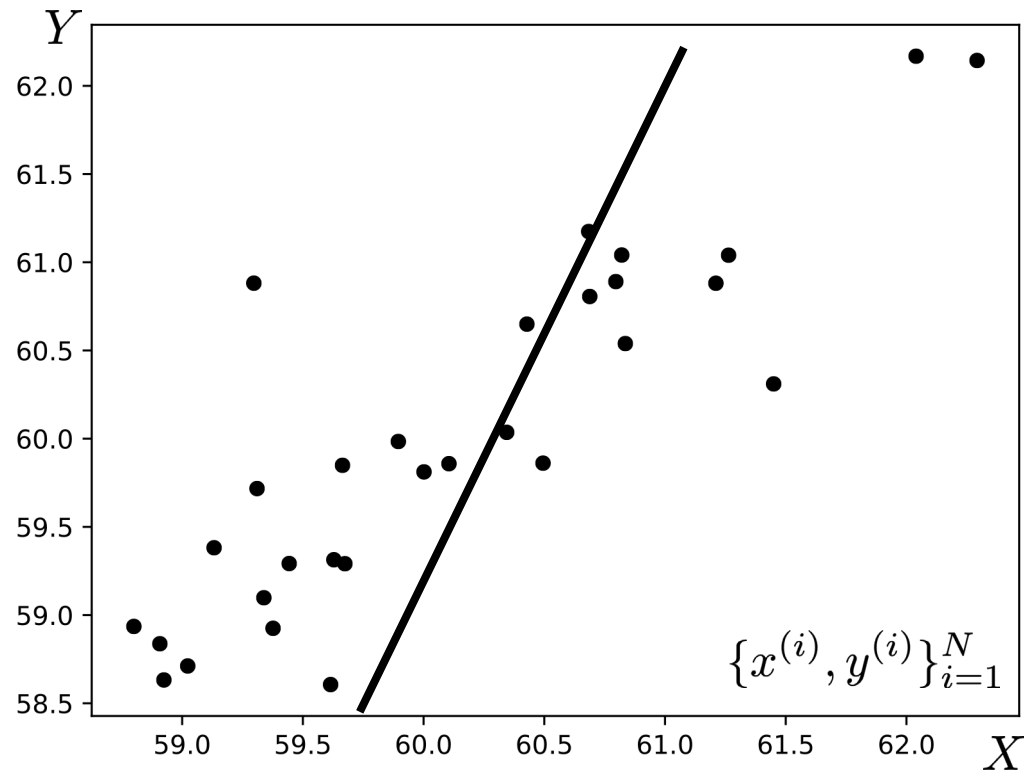


$$f_{\theta}(x) = \theta_1 x + \theta_0$$

Hypothesis space

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

Training data

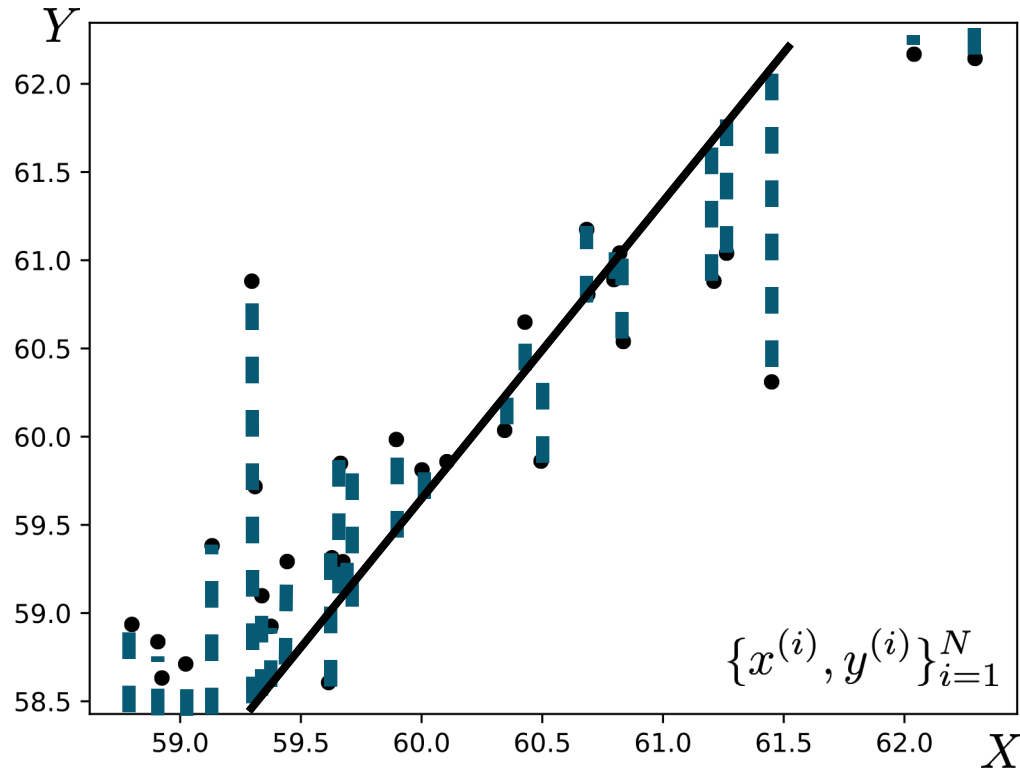


Search for the **parameters**, $\theta = \{\theta_0, \theta_1\}$, that best fit the data.

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

Best fit in what sense?

Training data



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

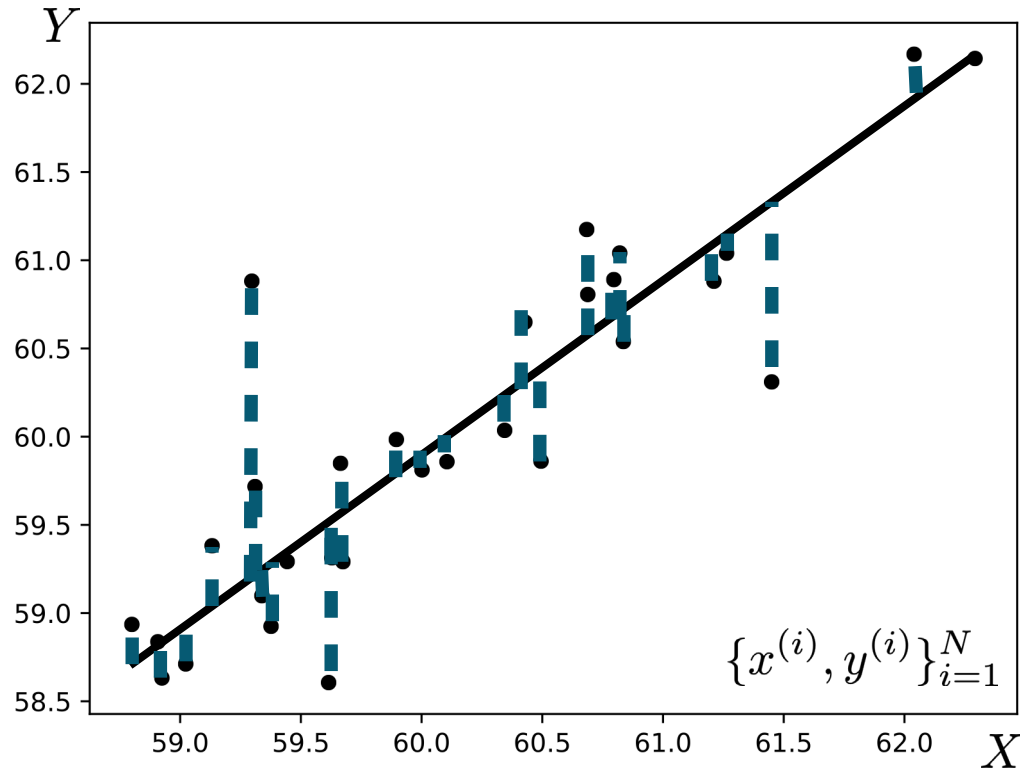
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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.

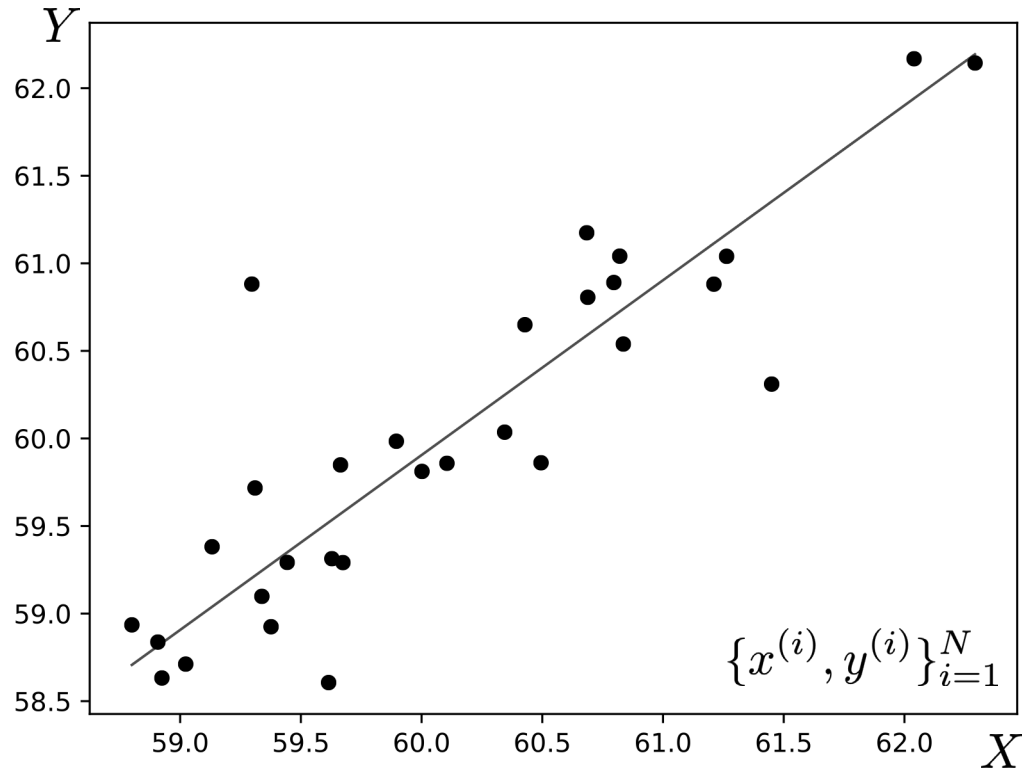
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

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:

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

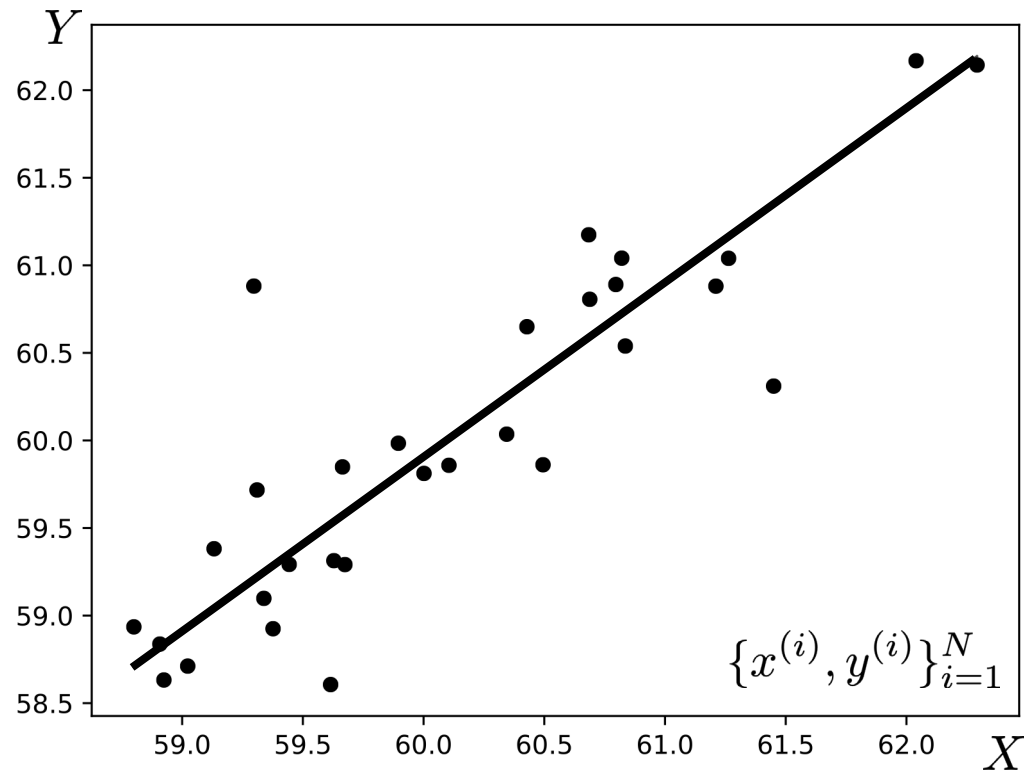
Training data



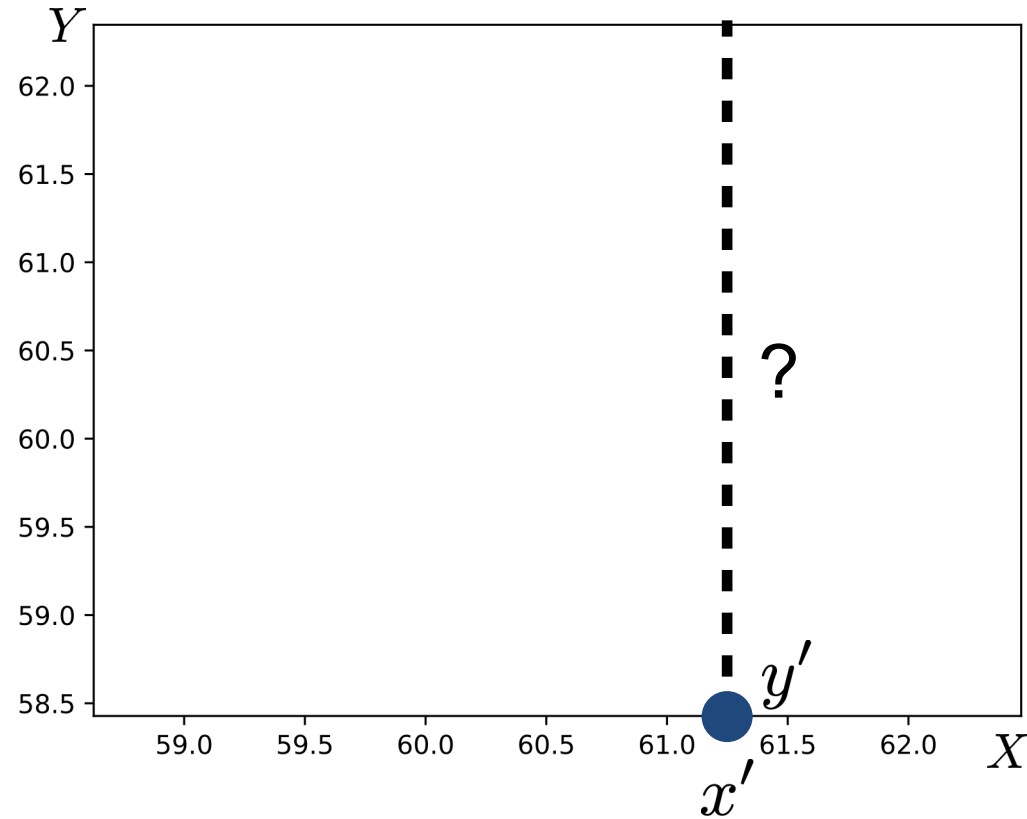
Complete learning problem:

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

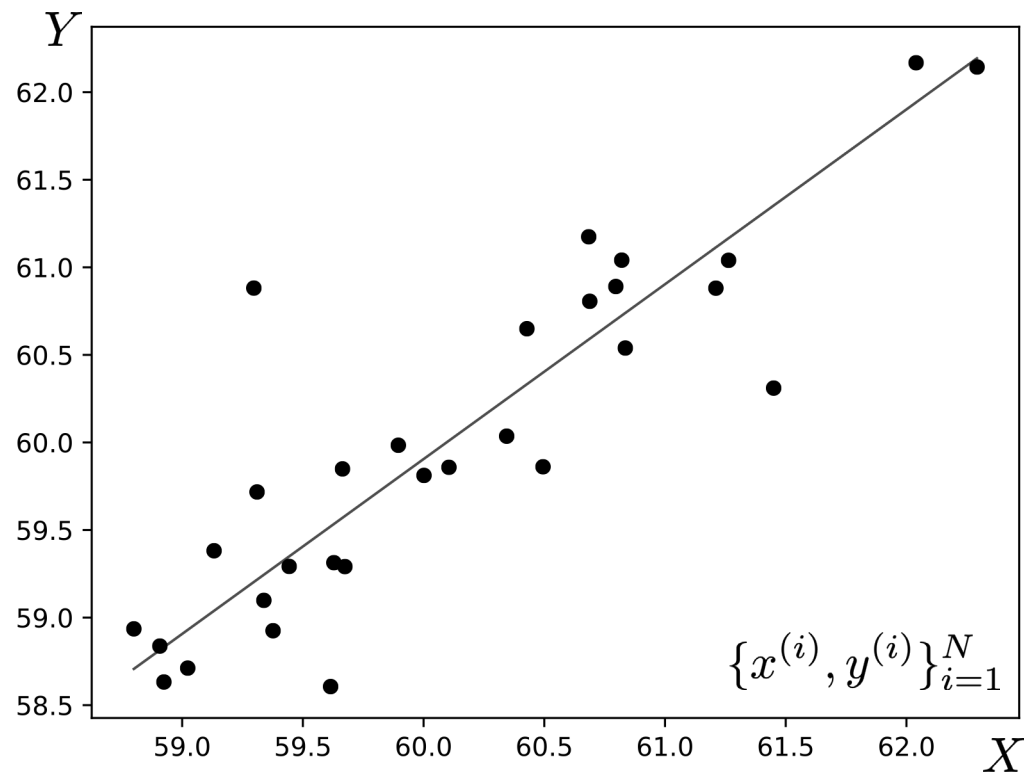
Training data



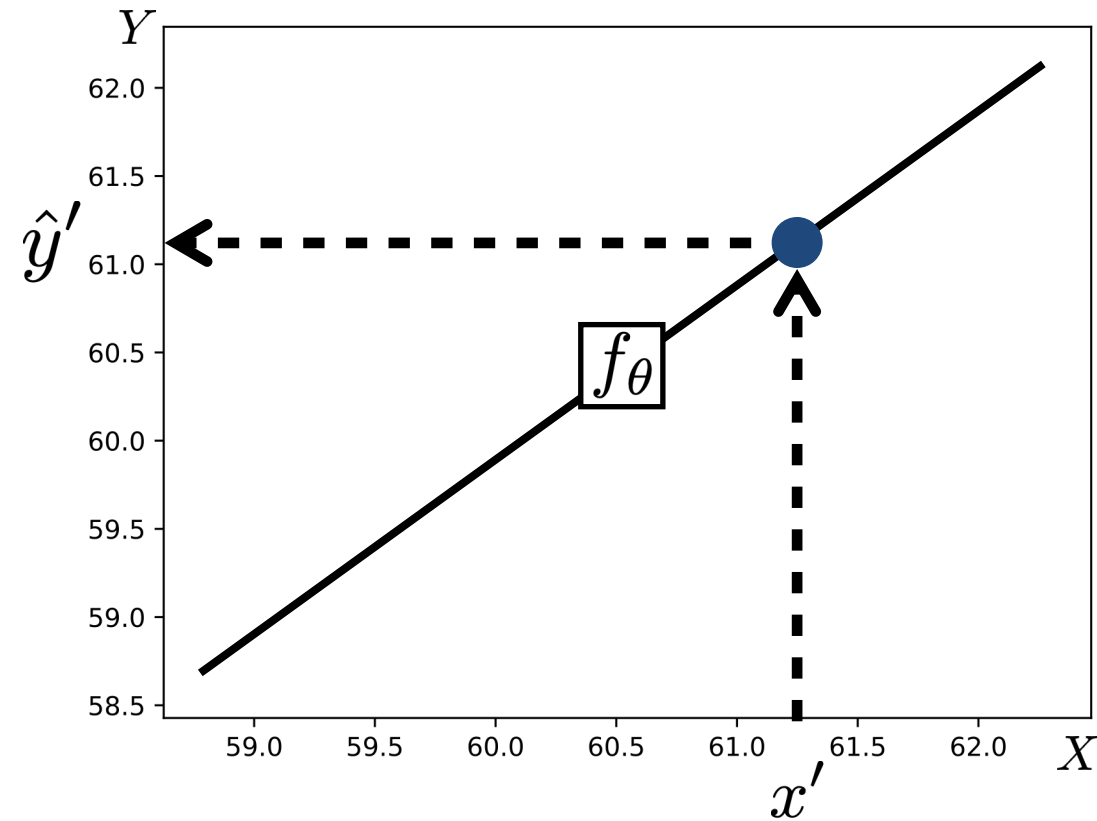
Test query



Training data



Test query

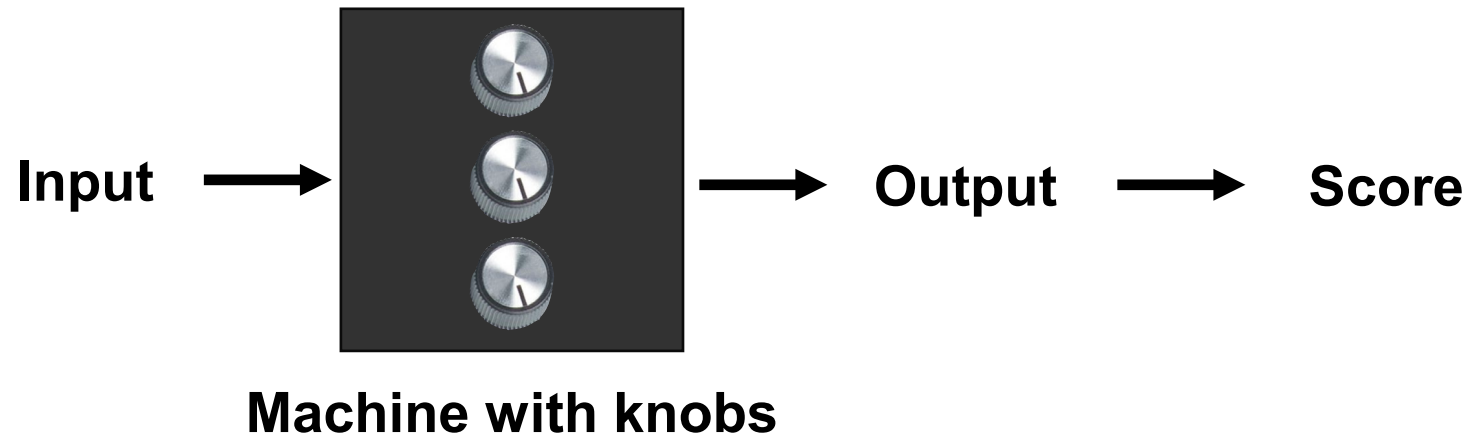


$$x' \dashrightarrow \boxed{f_\theta} \dashrightarrow \hat{y}'$$

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

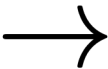
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x^{(i)}) - y^{(i)})^2$$

Use an **optimizer**!

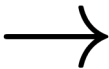


Learning

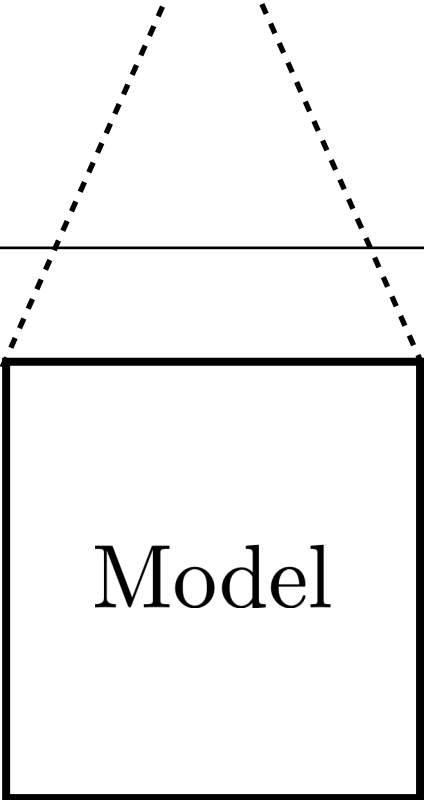
Data



Learner

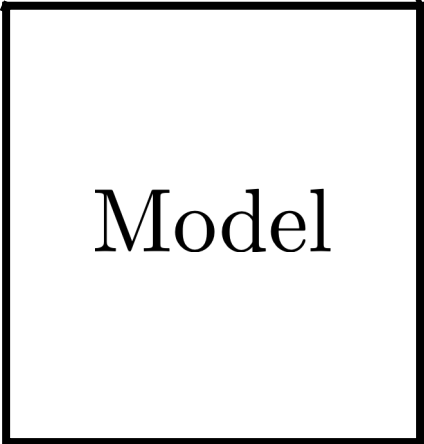
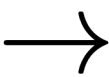


Model

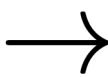


Inference

Input



Model



Output

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

In the linear case:

Learning problem

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

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

$$\mathbf{X} = \begin{pmatrix} x^{(1)} & 1 \\ x^{(2)} & 1 \\ \vdots & \vdots \\ x^{(N)} & 1 \end{pmatrix} \quad \theta = (\theta_1 \quad \theta_0) \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{pmatrix}$$

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

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

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

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

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

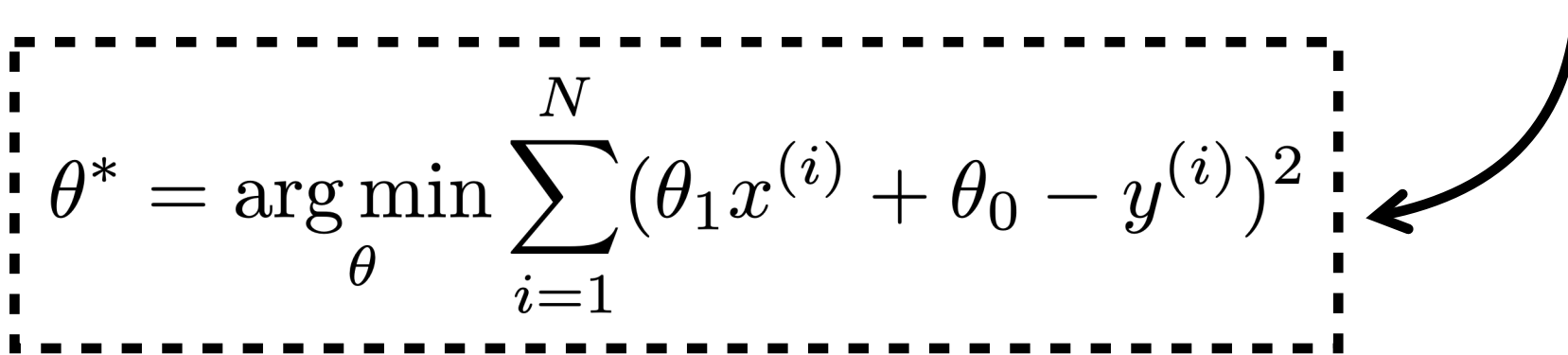
$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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)

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

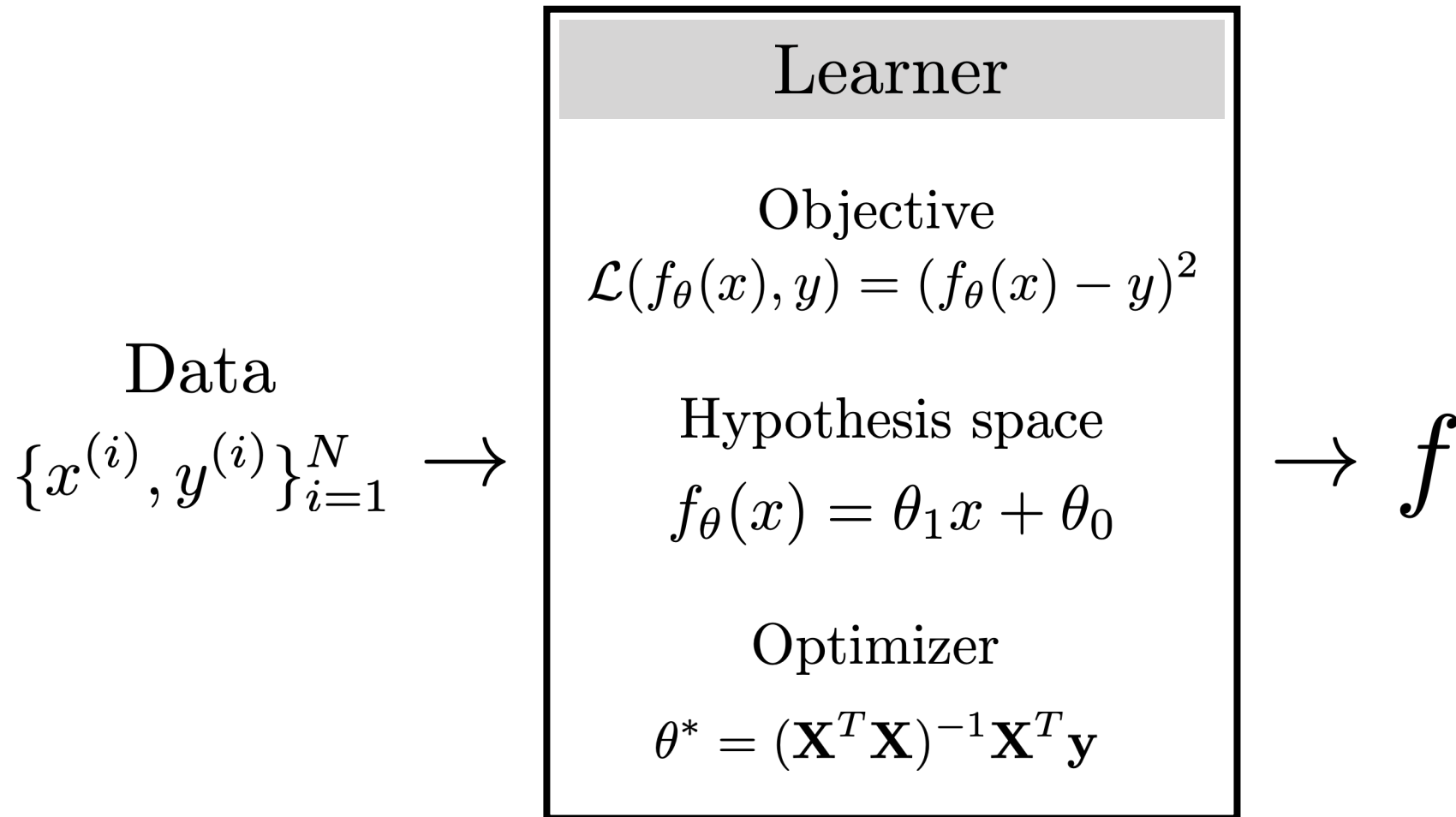
Objective function
(loss)

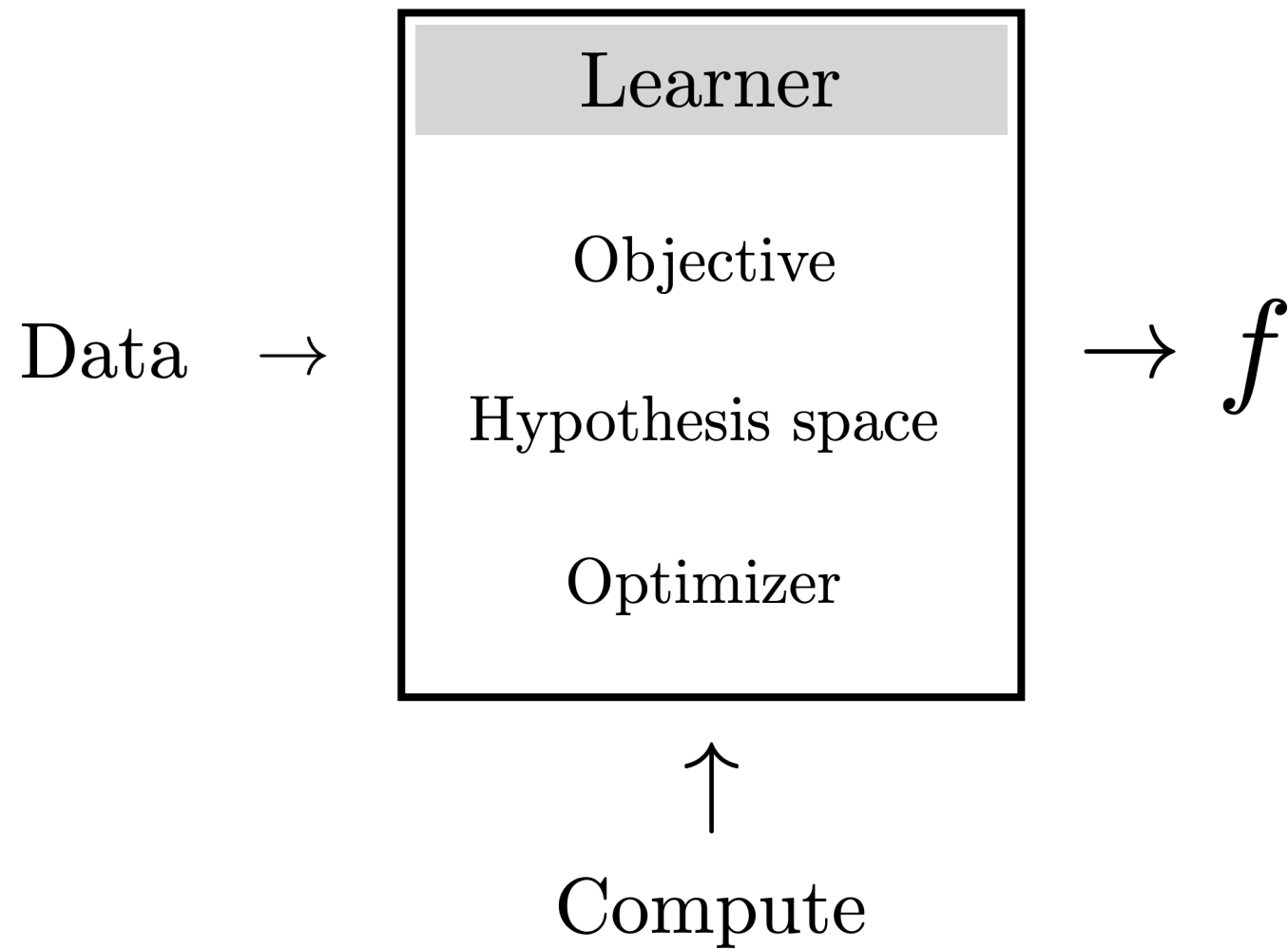
Hypothesis space

Training data

The diagram illustrates the Empirical Risk Minimization equation. The equation is $f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$. Annotations include: an arrow pointing from 'Objective function (loss)' to the loss function \mathcal{L} ; an arrow pointing from 'Hypothesis space' to the minimization set $f \in \mathcal{F}$; and two arrows pointing from 'Training data' to the input $\mathbf{x}^{(i)}$ and output $\mathbf{y}^{(i)}$ terms within the summation.

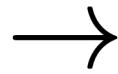
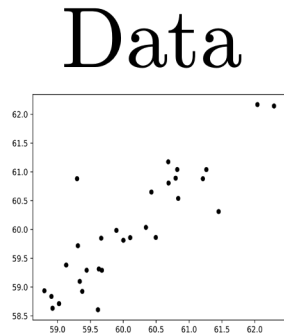
Case study #1: Linear least squares



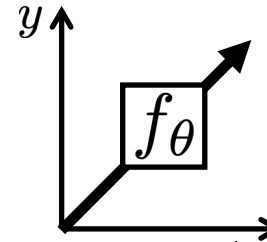
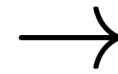


Example 1: Linear least squares

Training

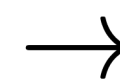
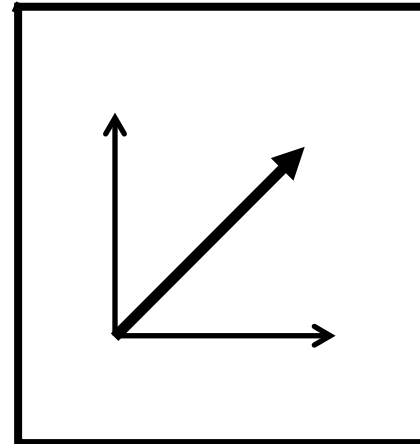
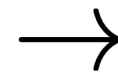


Learner



Testing

Input

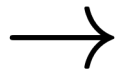
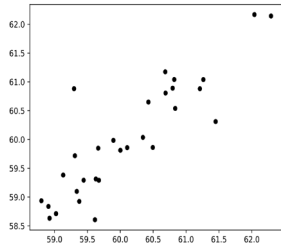


Output

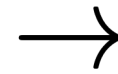
Example 2: Program Induction

Training

Data



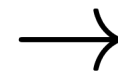
Learner



```
def predict(x):  
    y = 0.8*x + 2  
    return y
```

Testing

Input

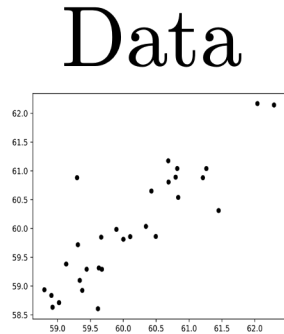


```
def predict(x):  
    y = 0.8*x + 2  
    return y
```

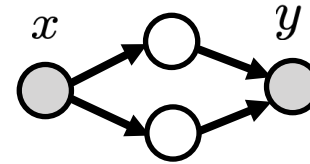
→ Output

Example 3: “Deep” Learning (with Neural Networks)

Training

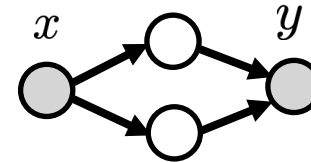


Learner



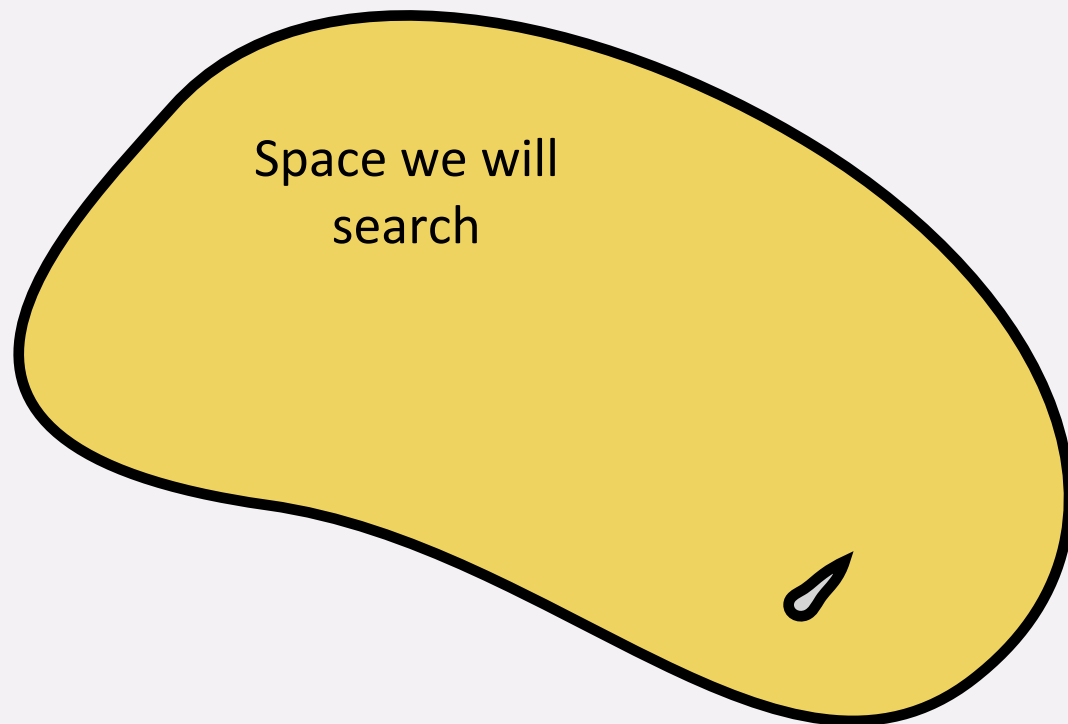
Testing

Input



Output

Space of all functions



Hypothesis space (haystack)



True solution (needle)

Space of all functions



Hypothesis space (haystack)



True solution (needle)

Space we will
search



Linear functions

True solution is linear

Space of all functions



Hypothesis space (haystack)



True solution (needle)

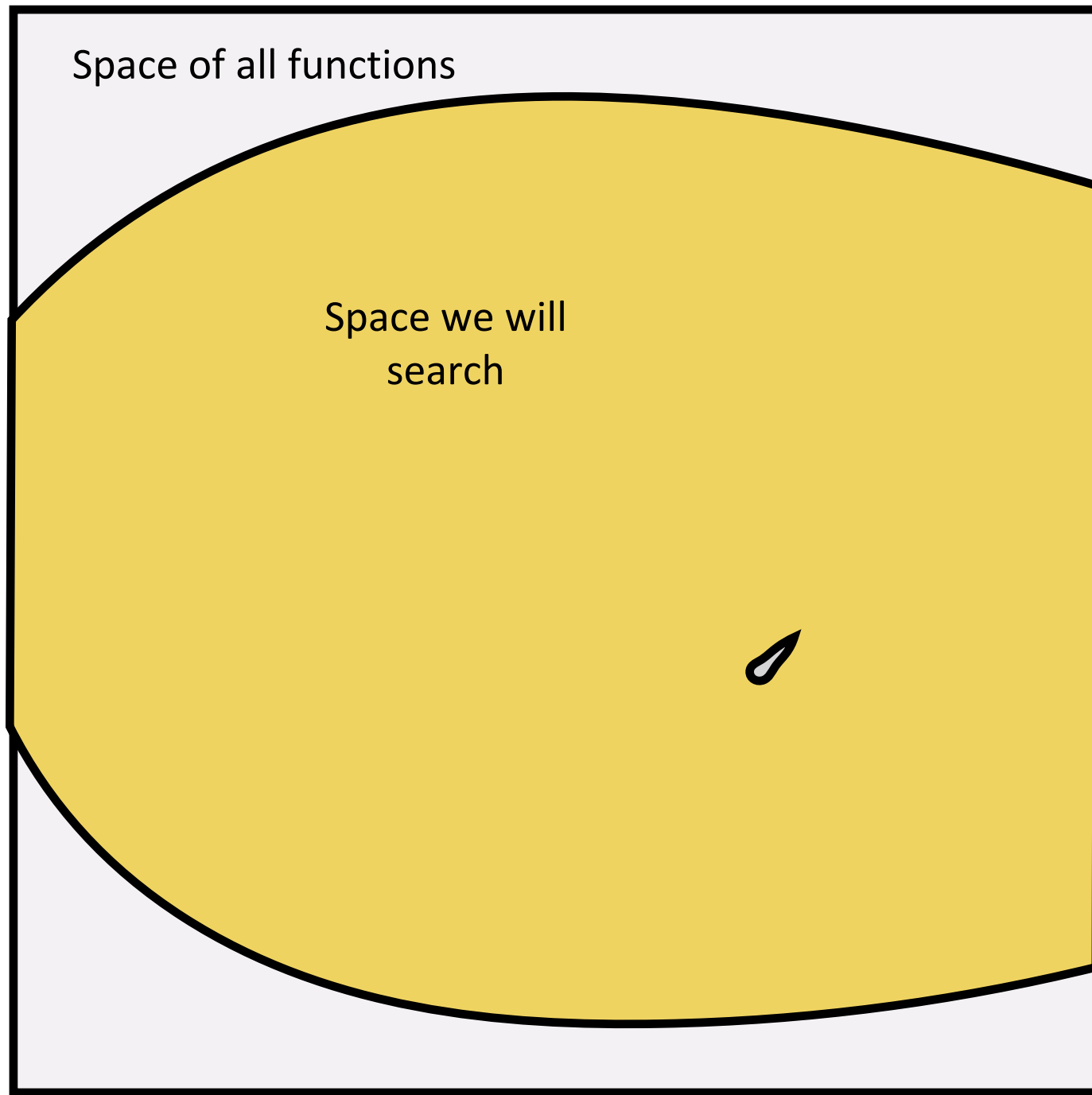


Space we will
search



Linear functions

True solution is nonlinear



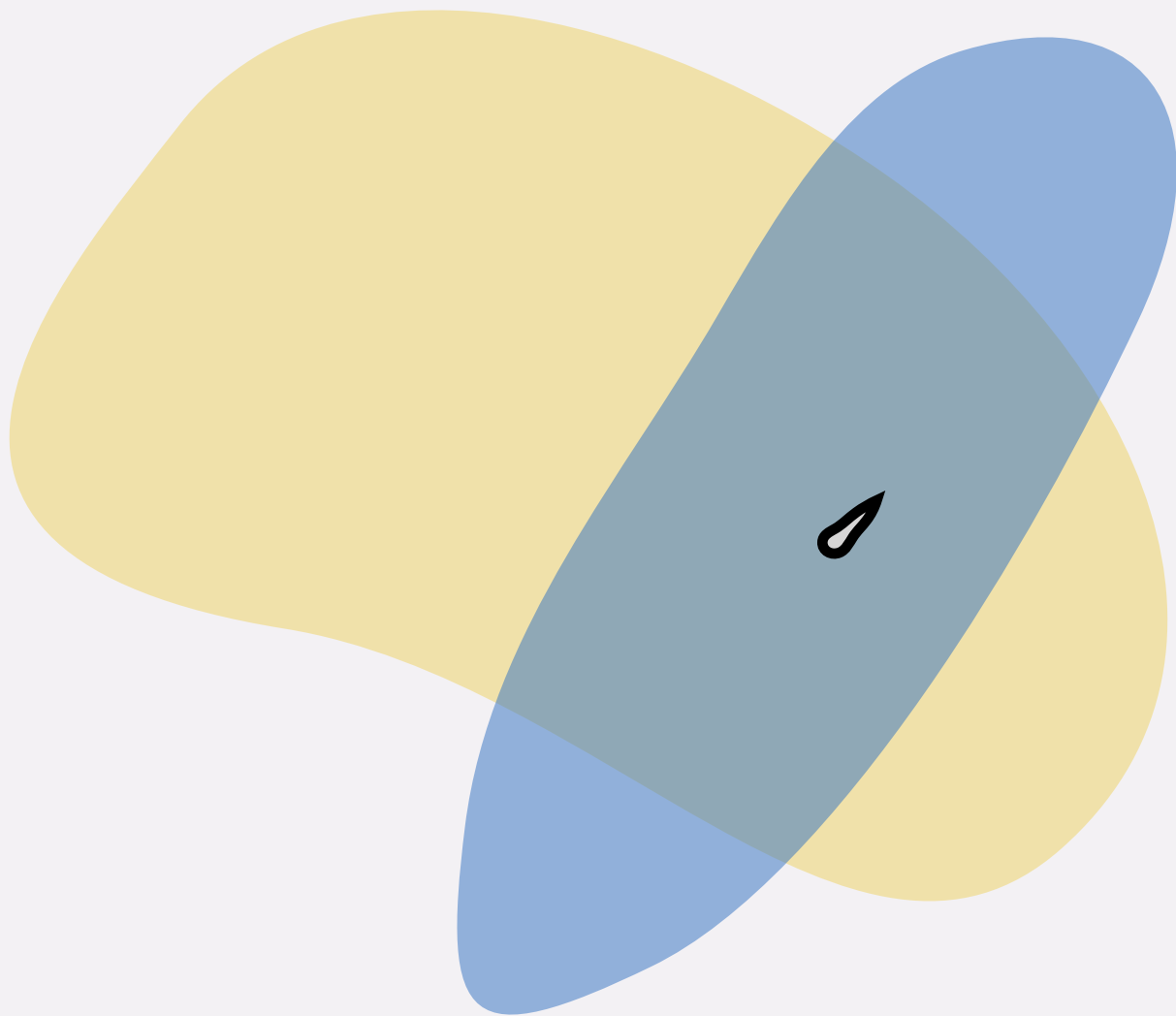
Hypothesis space (haystack)



True solution (needle)

Deep nets

Space of all functions



Hypothesis space (haystack)



True solution (needle)



Hypotheses consistent with data

Space of all functions



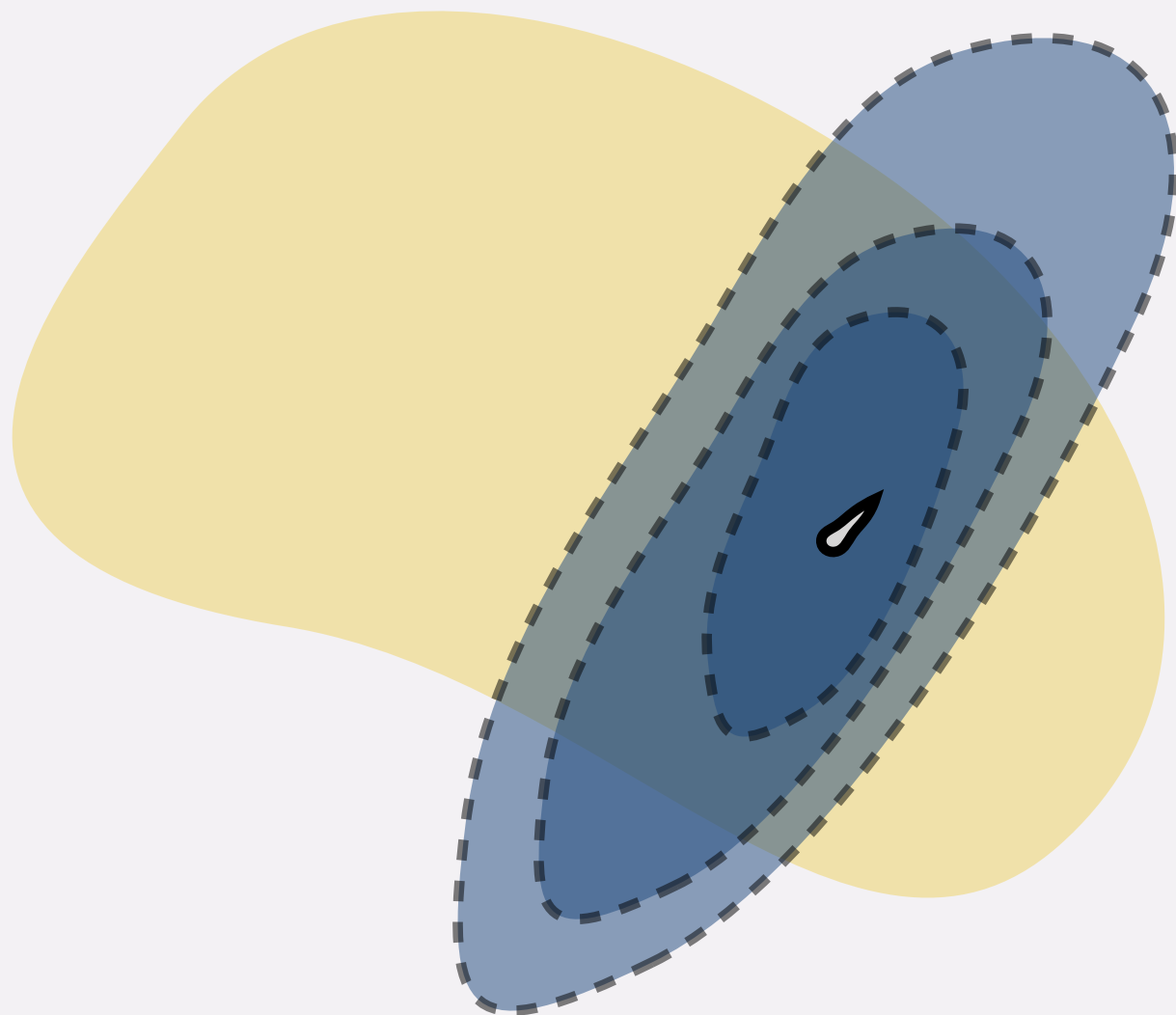
Hypothesis space (haystack)



True solution (needle)



Hypotheses consistent with data



What happens as we increase
the data?

Space of all functions



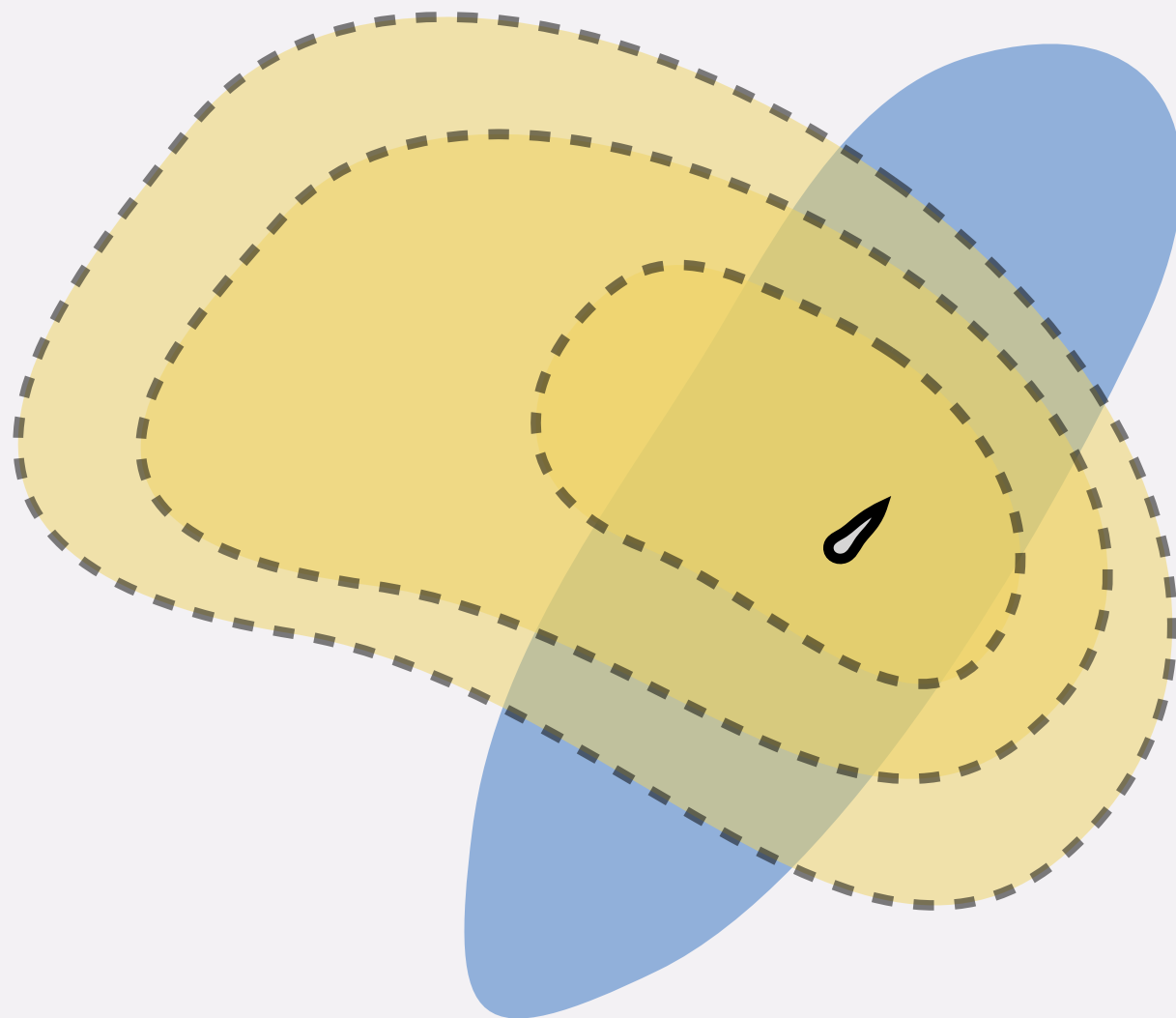
Hypothesis space (haystack)



True solution (needle)



Hypotheses consistent with data



What happens as we shrink the hypothesis space?

The essence of machine learning:

- A pattern exists
- We cannot pin down the pattern as an equation
- We need to approximate the pattern as a function of the input
 - Using a set of observations (data) to uncover an underlying process

Regression vs. Classification

- Regression tasks: predicting real-valued outputs $y \in \mathbb{R}$
- Classification tasks: predicting discrete-valued quantity
 y
 - Binary Classification $y \in \{-1, 1\}$
 - Multiclass Classification $y \in \{1, 2, \dots, k\}$

Learning for vision

Big 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?

Case study #2: Image classification

- 1. How do you represent the input and output?**
- 2. What is the objective?**
3. Assume hypothesis space is sufficiently expressive
4. Assume we optimize perfectly
5. Assume we train on exactly the data we care about

Image classification

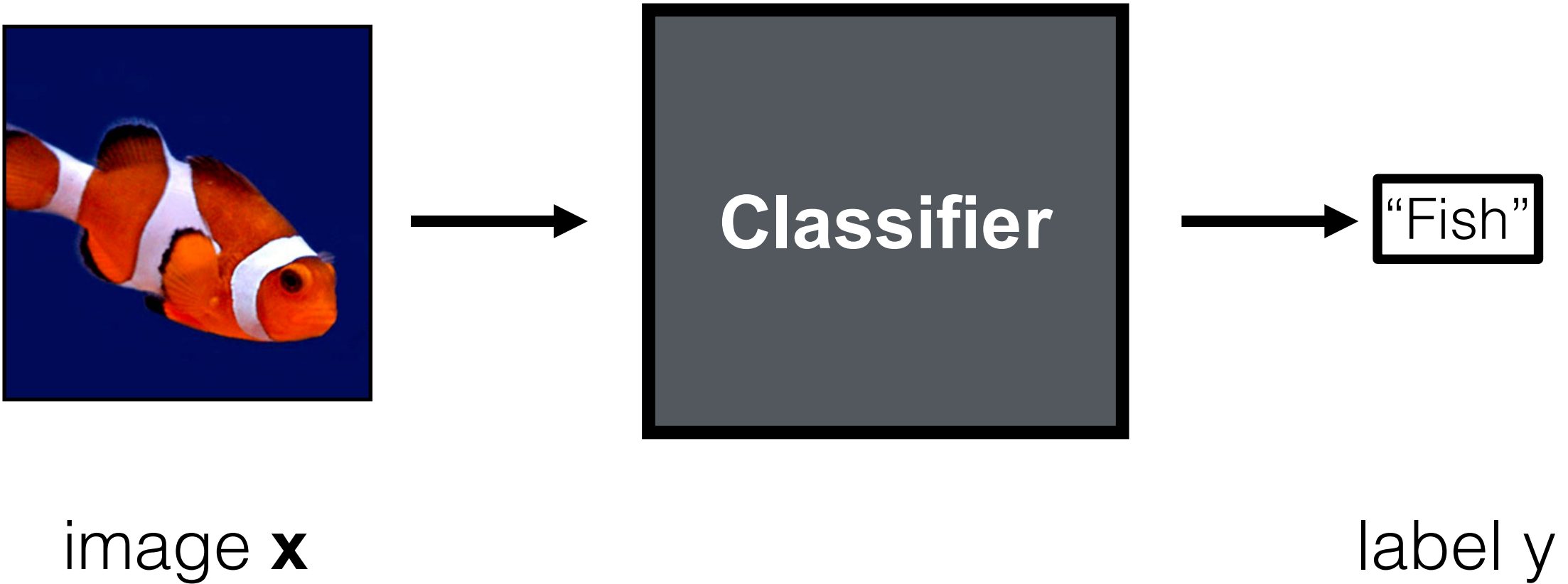


Image classification



image **x**



Classifier



"Fish"

label **y**

Image classification



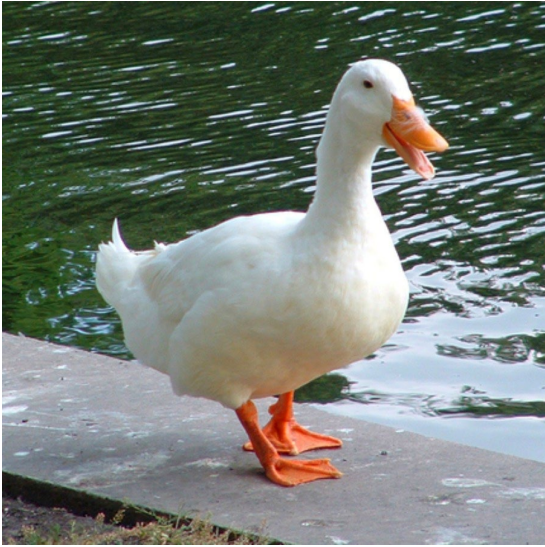
image **x**



"Fish"

label **y**

Image classification



⋮

image **x**



Classifier



“Duck”

label **y**

\mathbf{x}



Training data

\mathbf{x}

y



"Fish"

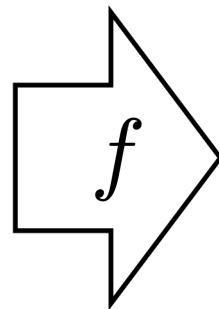


"Grizzly"



"Chameleon"

\vdots



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

Training data

\mathbf{x}

y



\vdots

Training data

\mathbf{x}

y

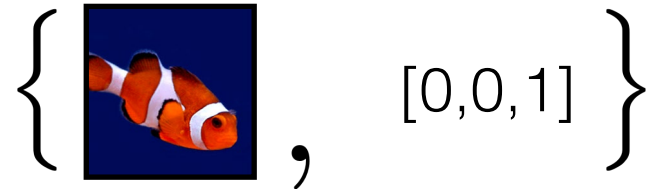


\vdots

Training data

\mathbf{x}

y



\vdots

What should the loss be?

0-1 loss (number of misclassifications)

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

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

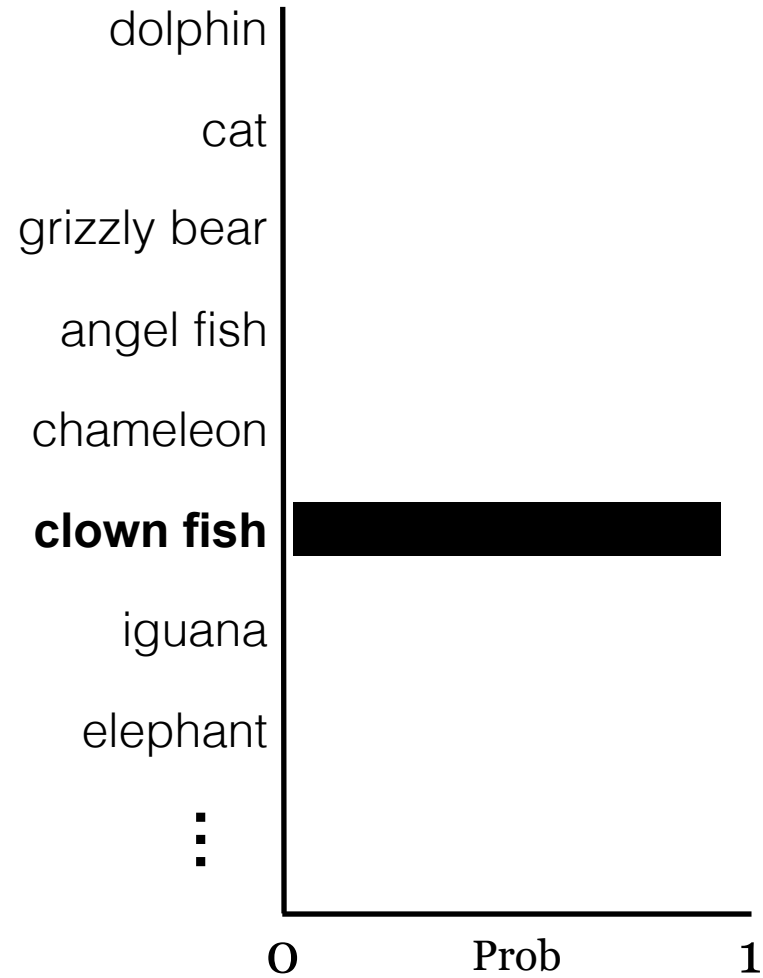
x



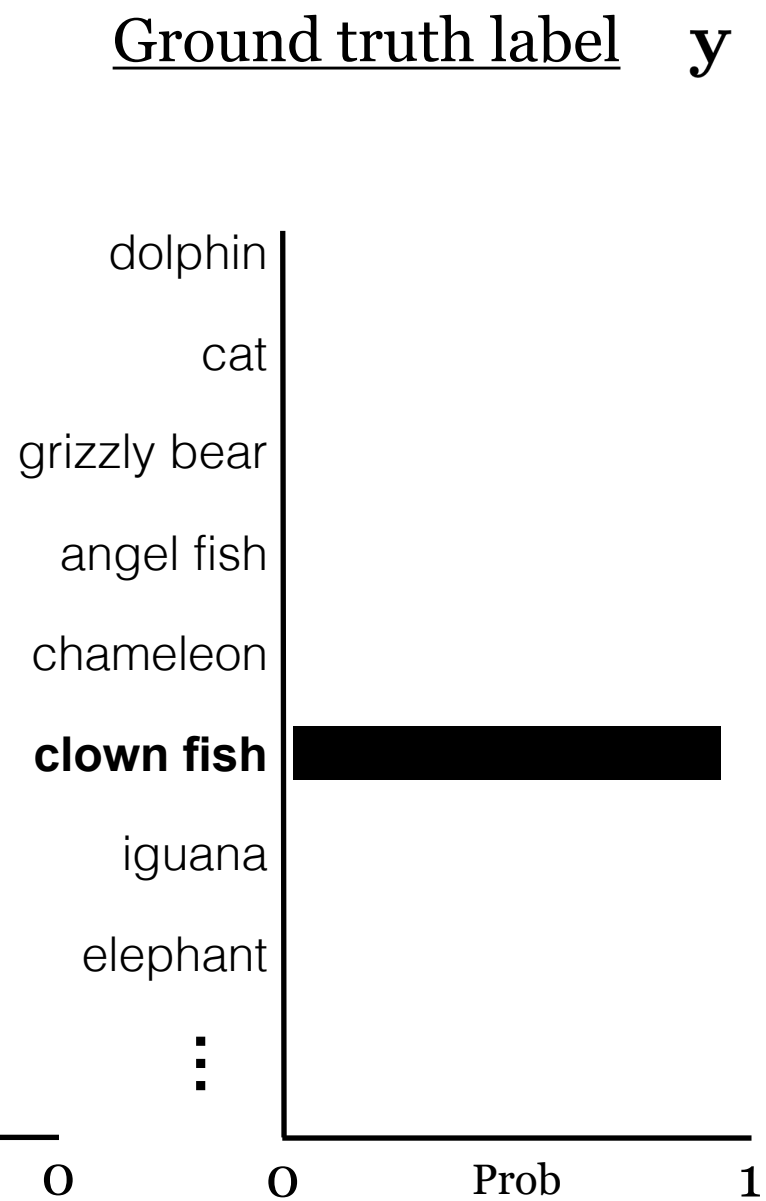
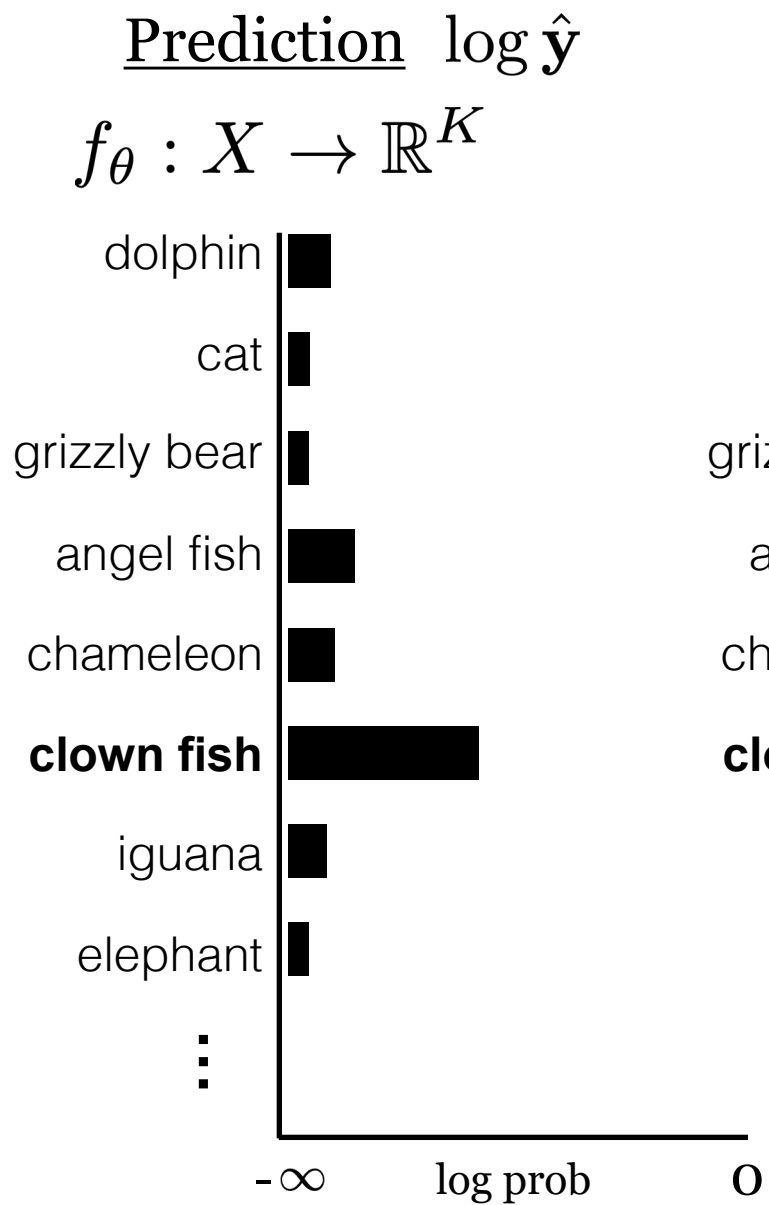
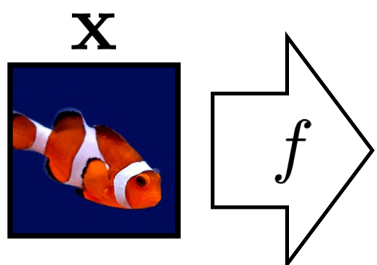
$[0,0,0,0,0,1,0,0,\dots]$

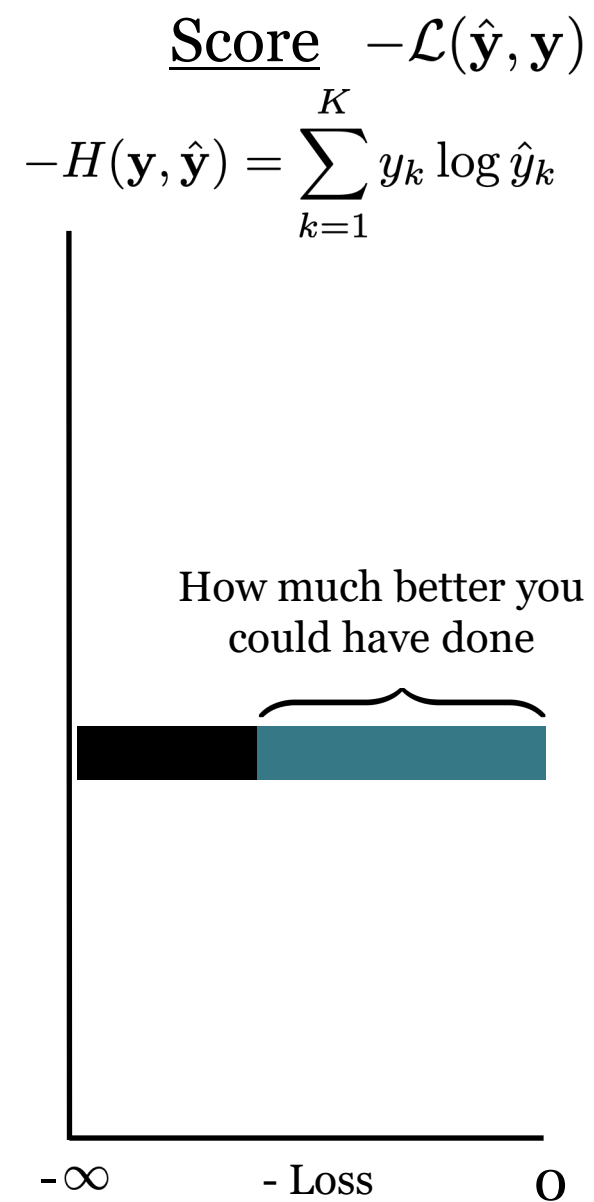
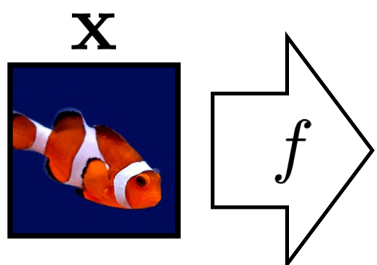


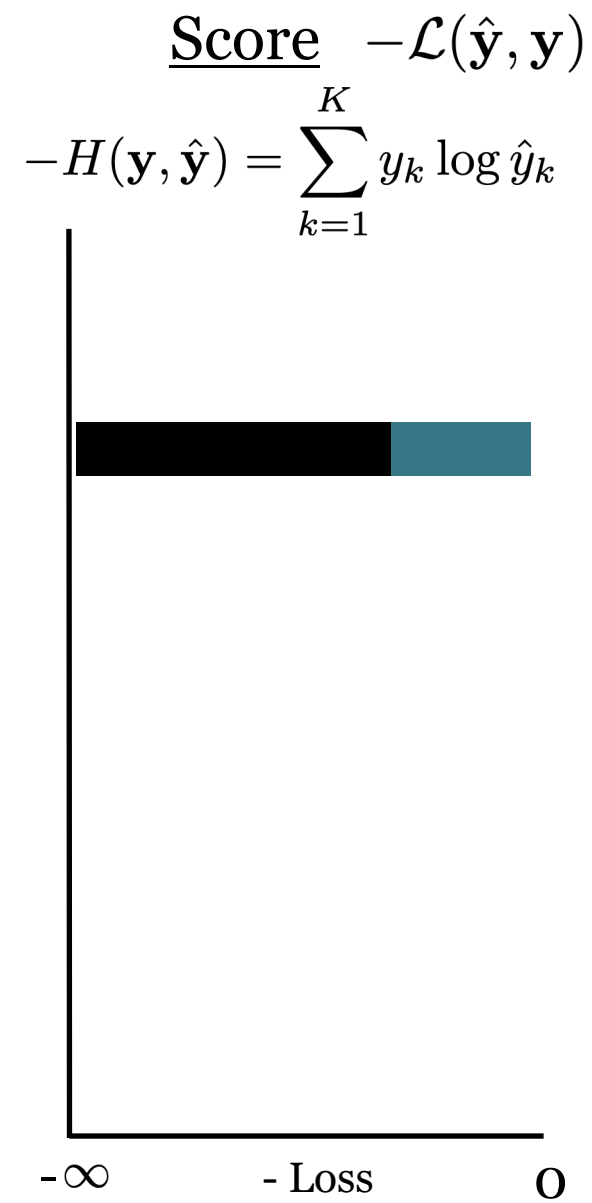
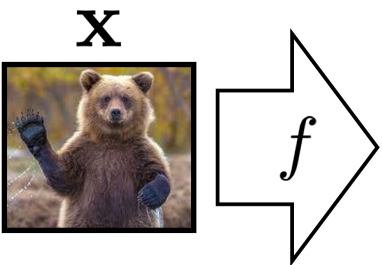
Ground truth label **y**

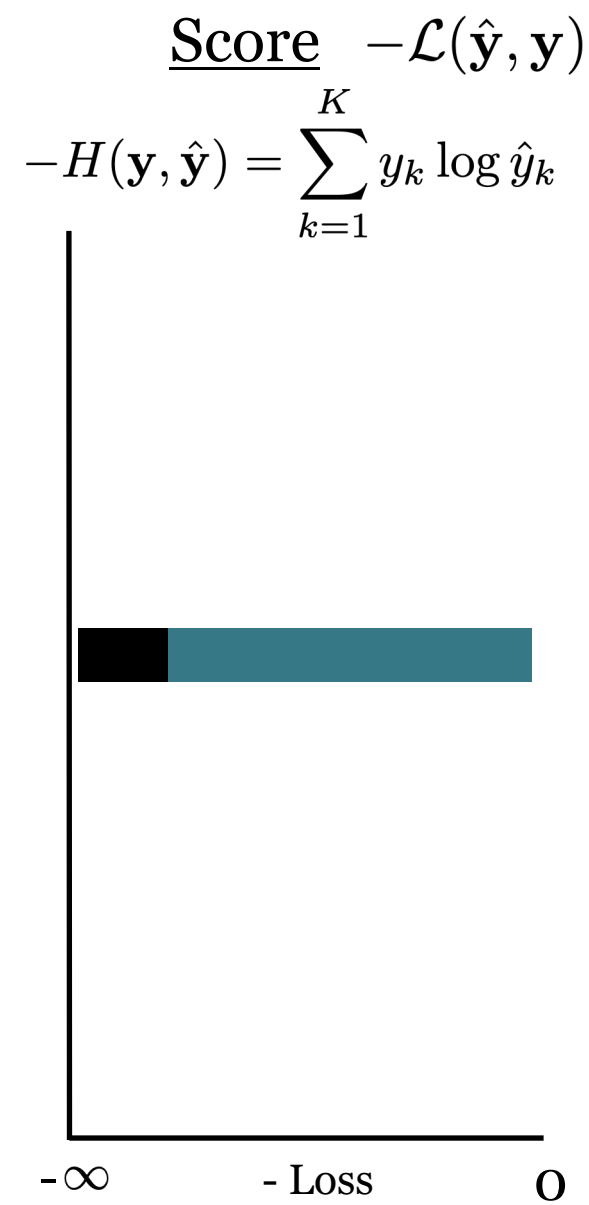
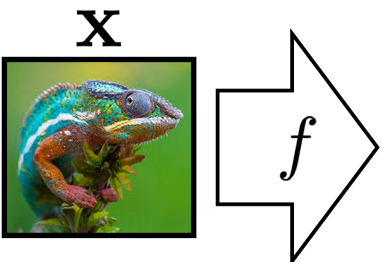


∞









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!**

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

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



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

Probabilistic interpretation:

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

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

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

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

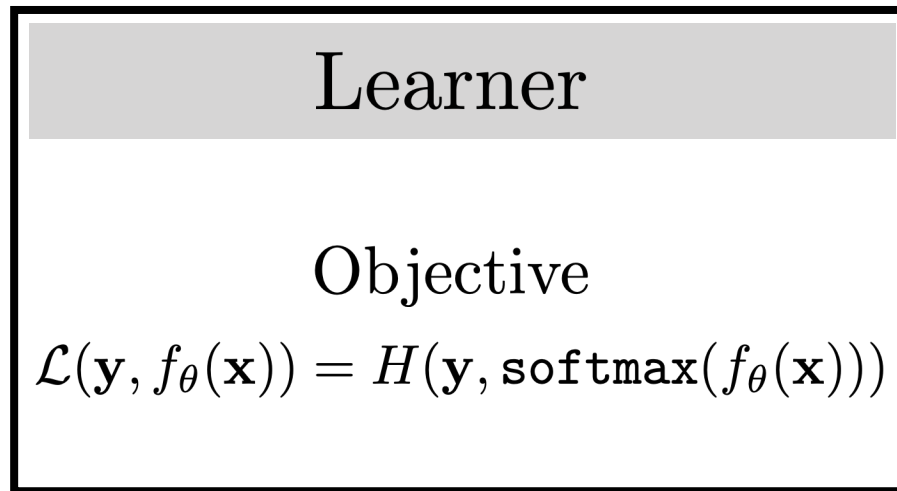
$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

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

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

Data

$$\{x^{(i)}, y^{(i)}\}_{i=1}^N \rightarrow$$



$$\rightarrow f$$

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.)

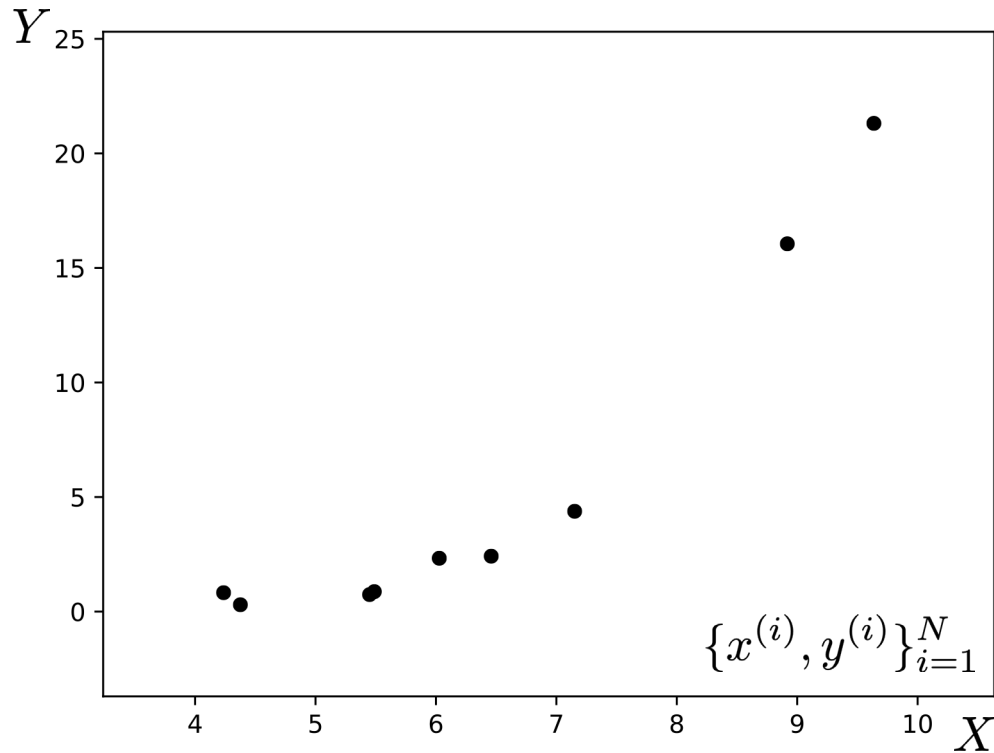
Recap:

Linear Regression

(f_{θ} is a linear function)

Linear regression

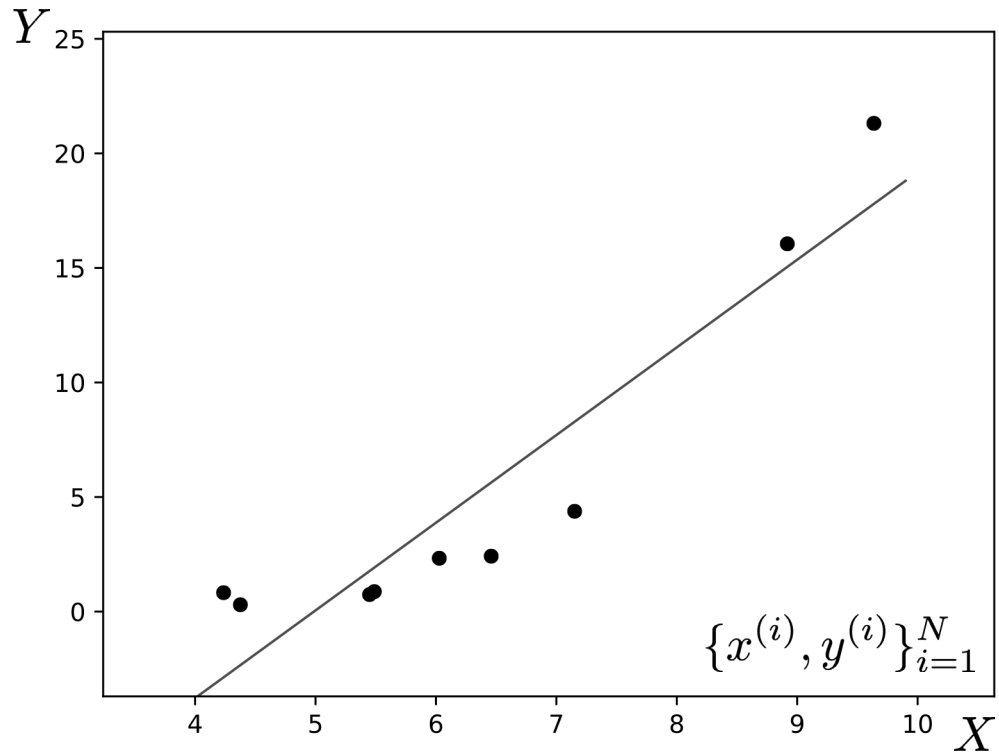
Training data



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

Linear regression

Training data



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

Linear Regression

(f_{θ} is a linear function)

Linear Regression

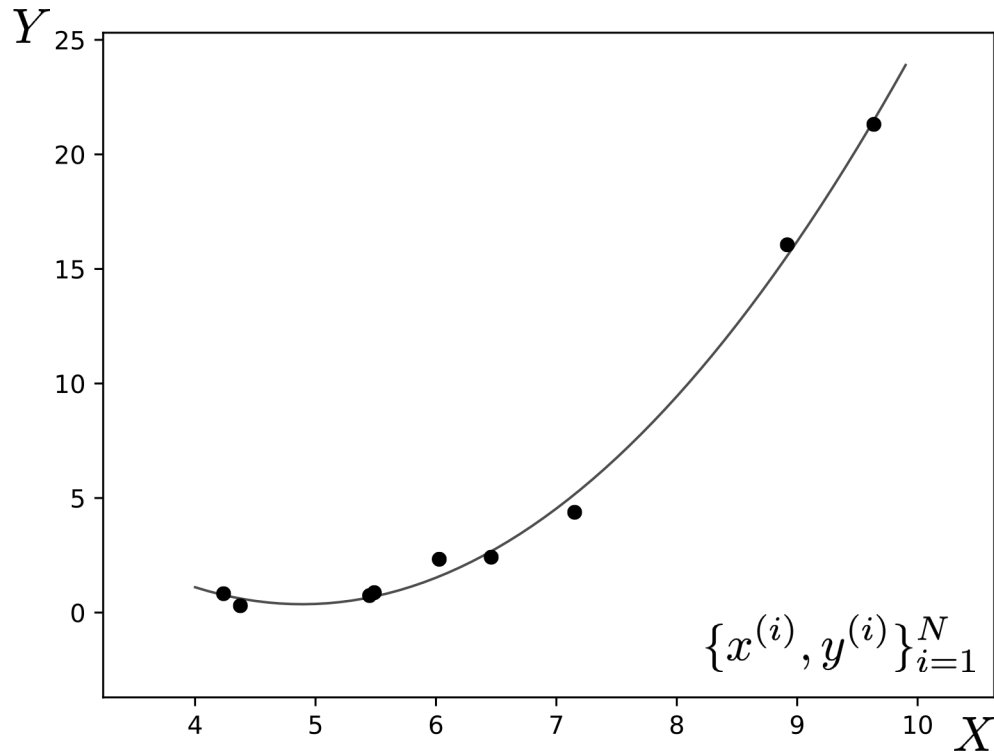
(f_{θ} is a linear function)

Polynomial Regression

(f_{θ} is a polynomial function)

Polynomial regression

Training data

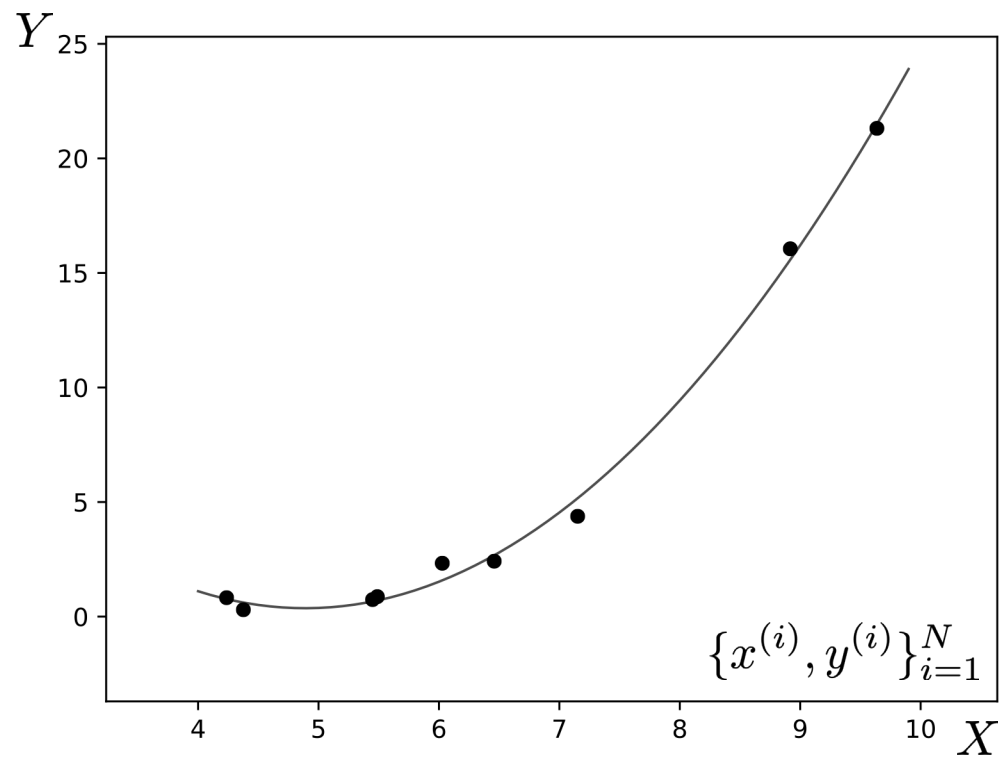


$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

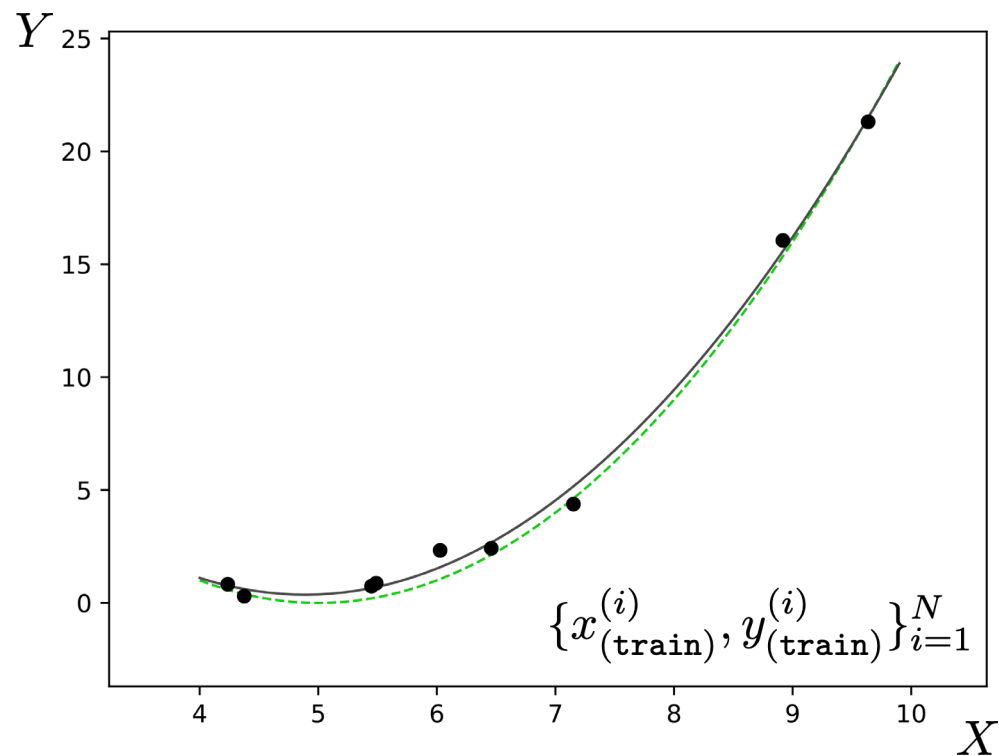
$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

K-th degree polynomial regression

Training data



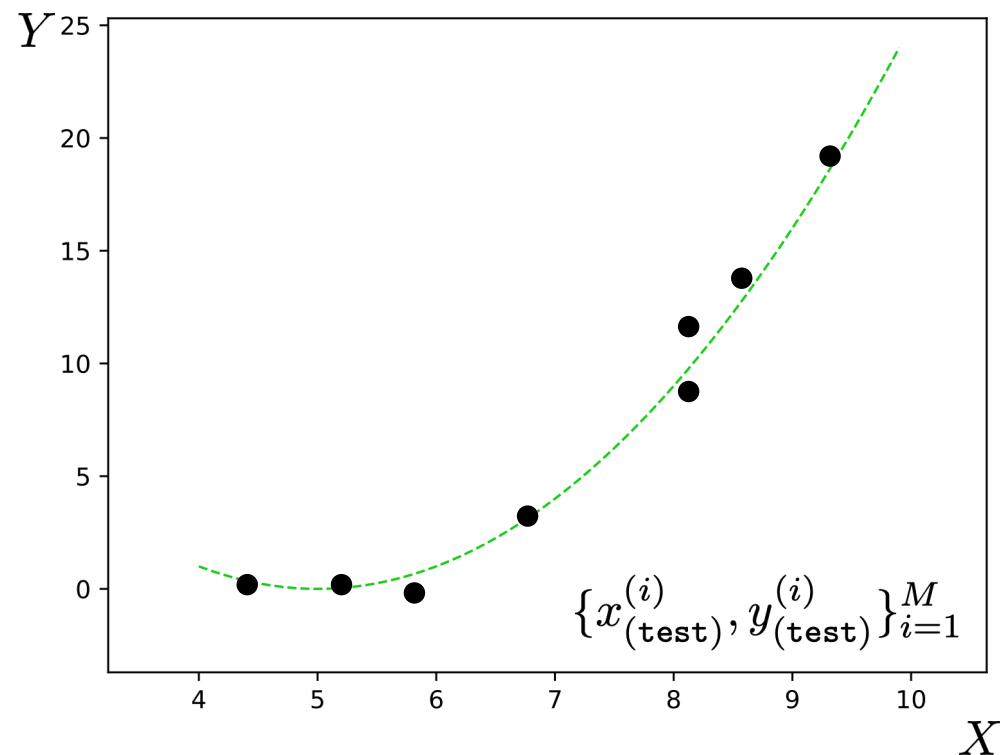
Training data



Training objective:

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

Test data

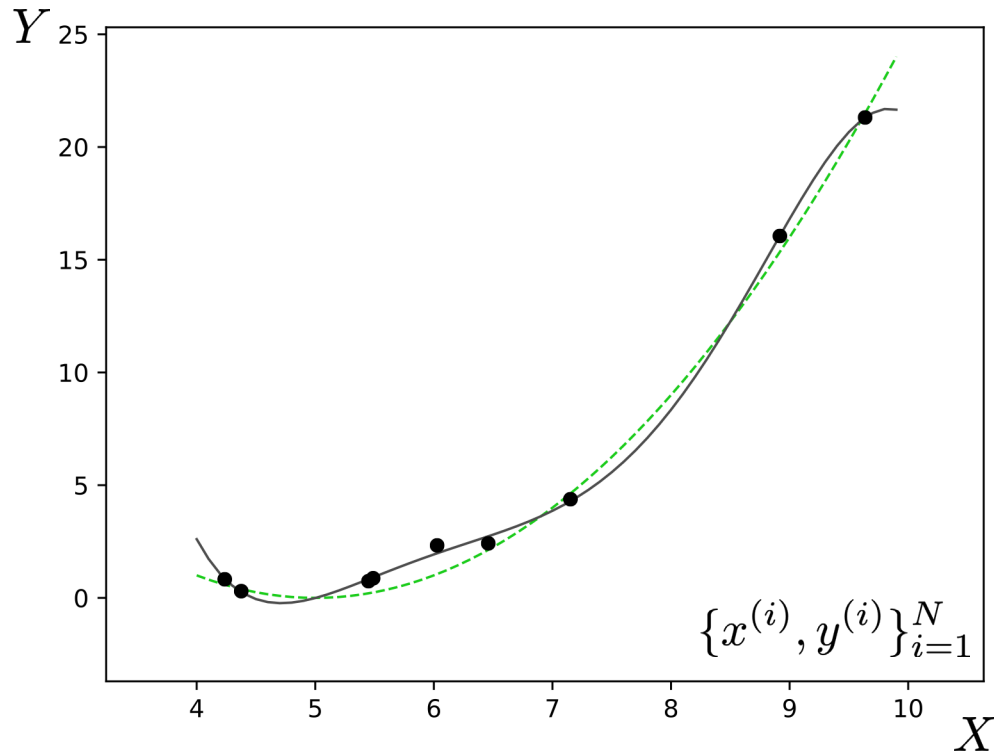


Test time evaluation:

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

What happens as we add more basis functions?

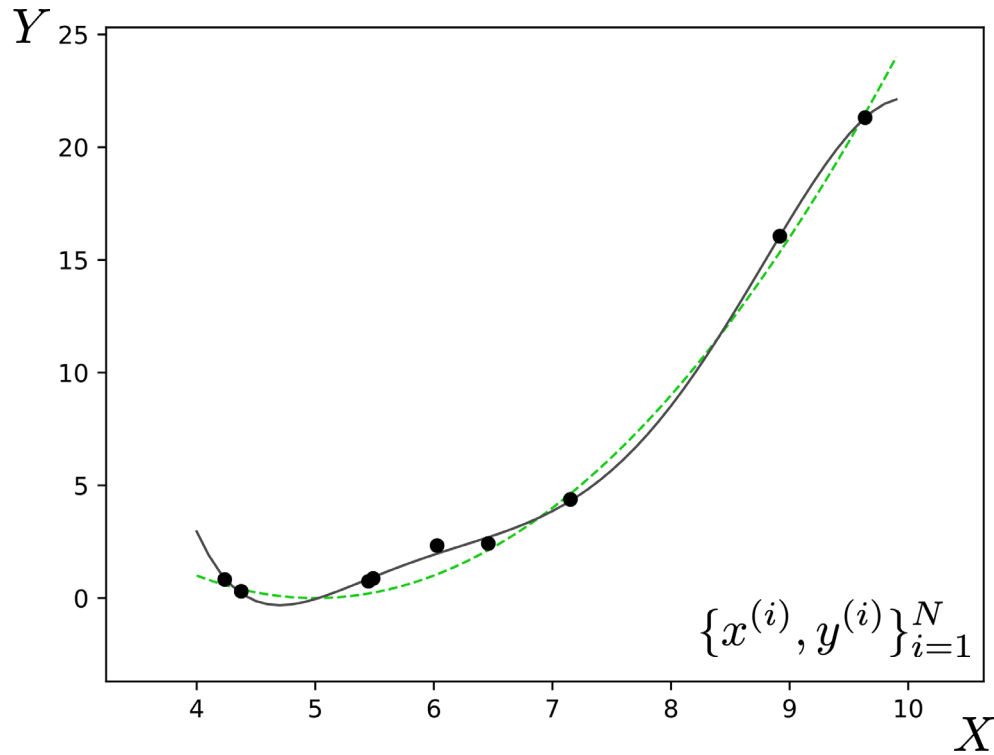
K = 5



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

What happens as we add more basis functions?

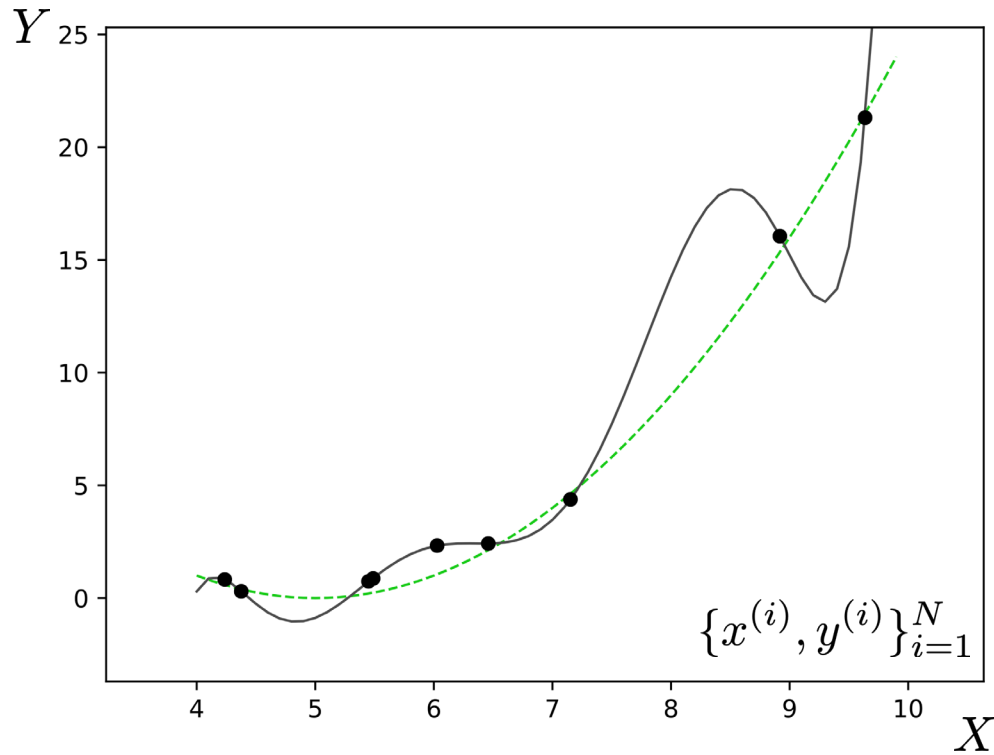
K = 6



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

What happens as we add more basis functions?

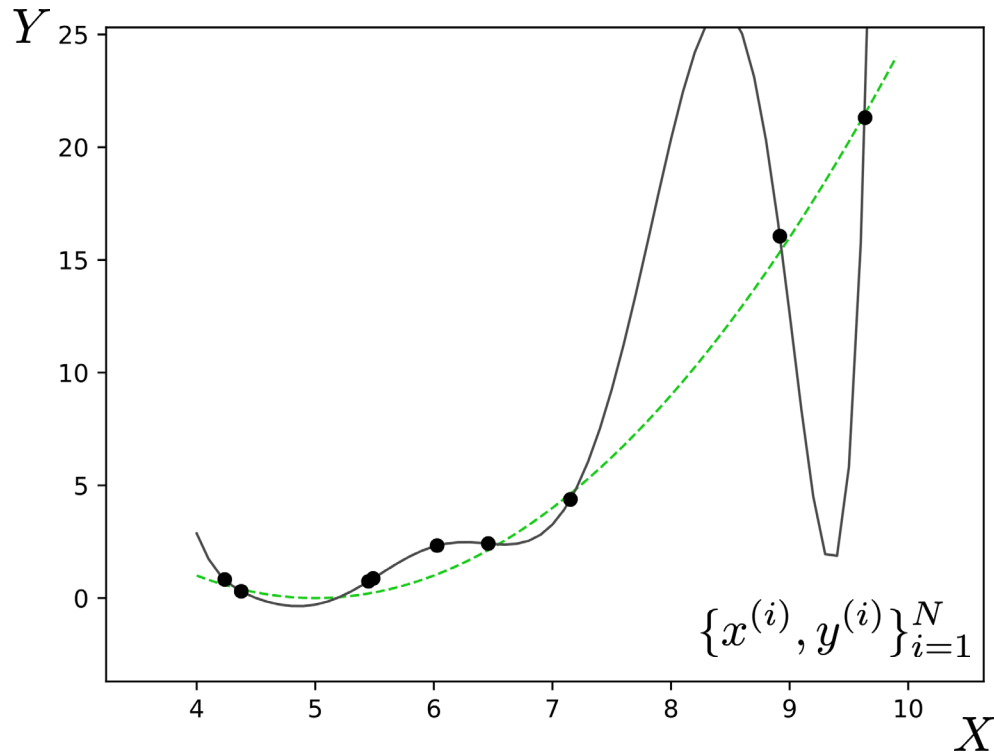
K = 7



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

What happens as we add more basis functions?

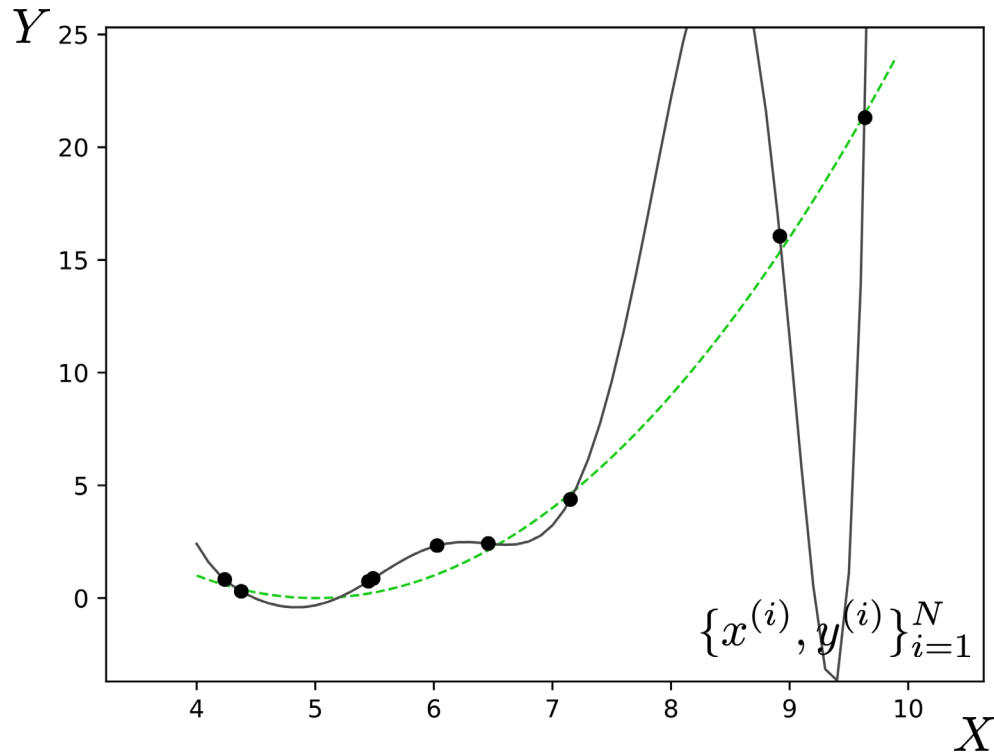
K = 8



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

What happens as we add more basis functions?

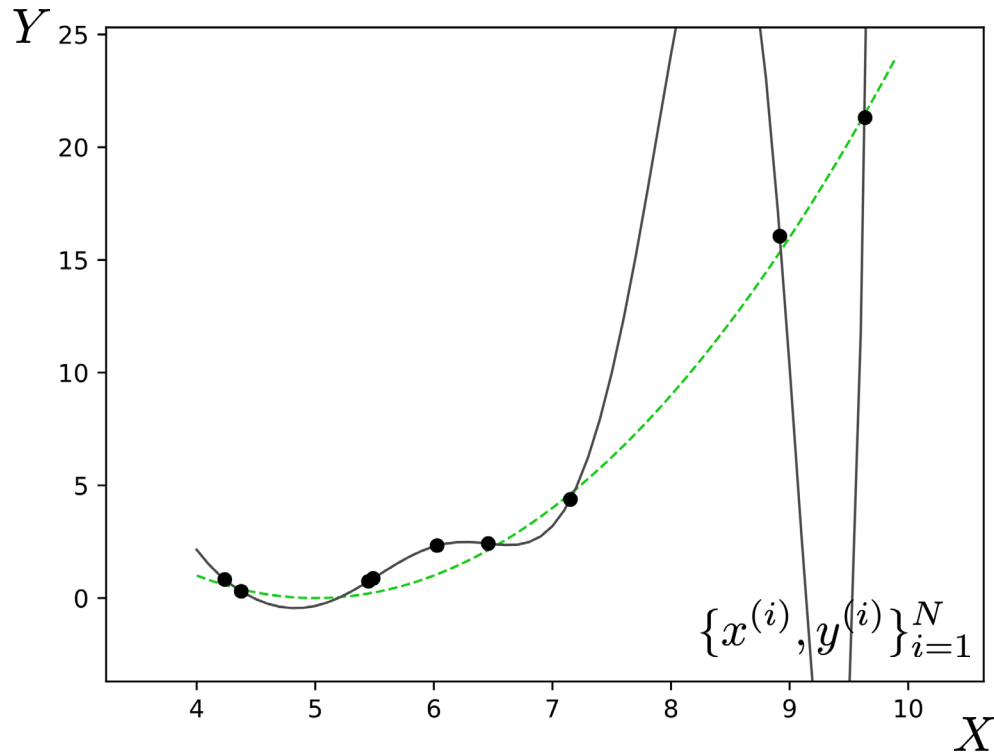
K = 9



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

What happens as we add more basis functions?

K = 10

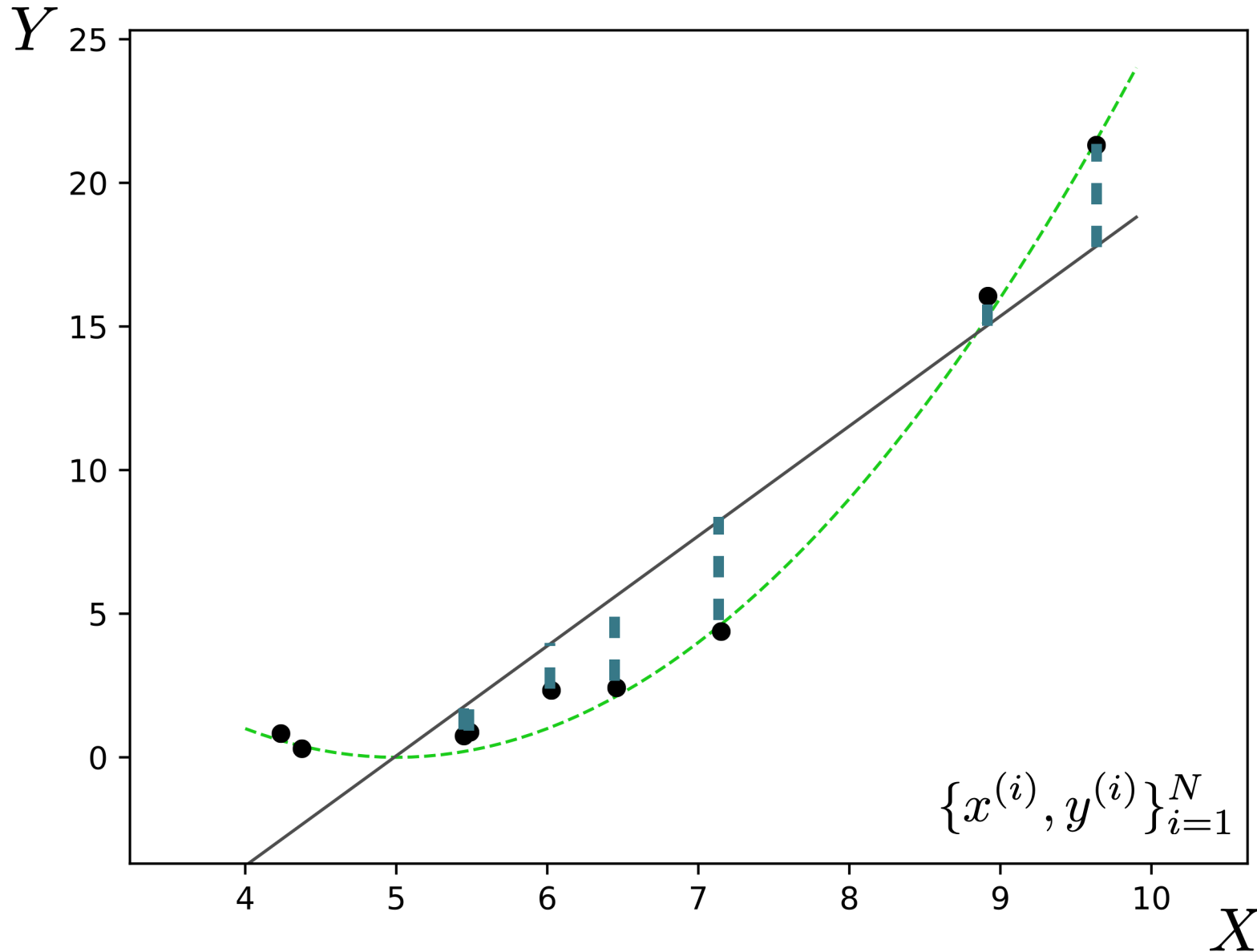


$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

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.

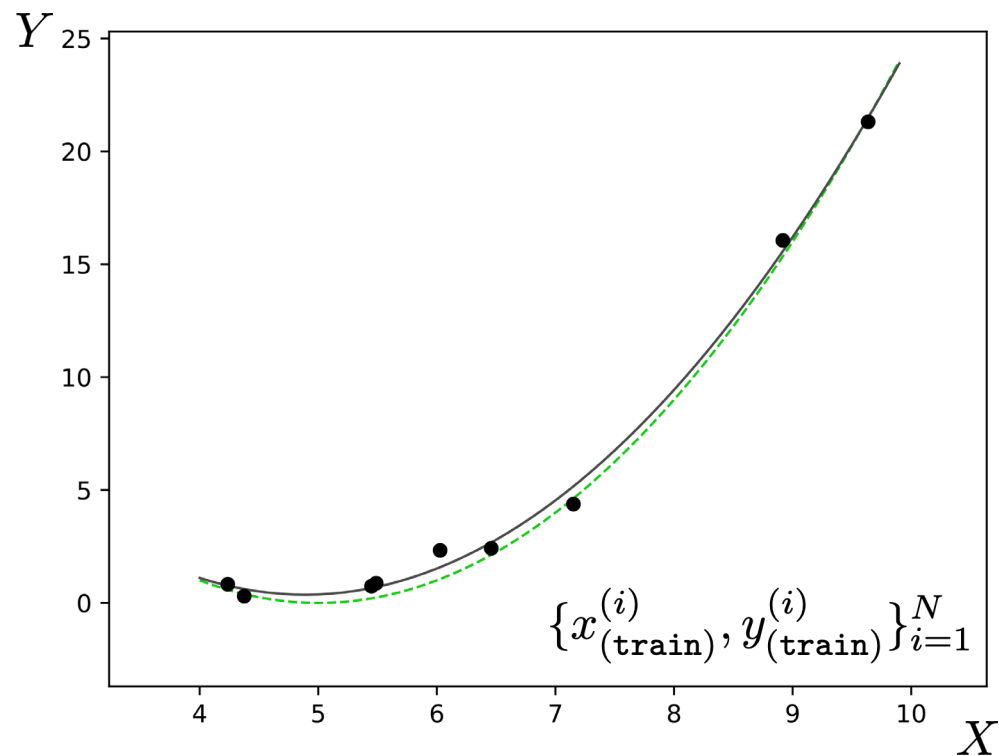
$$K = 1$$



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**.

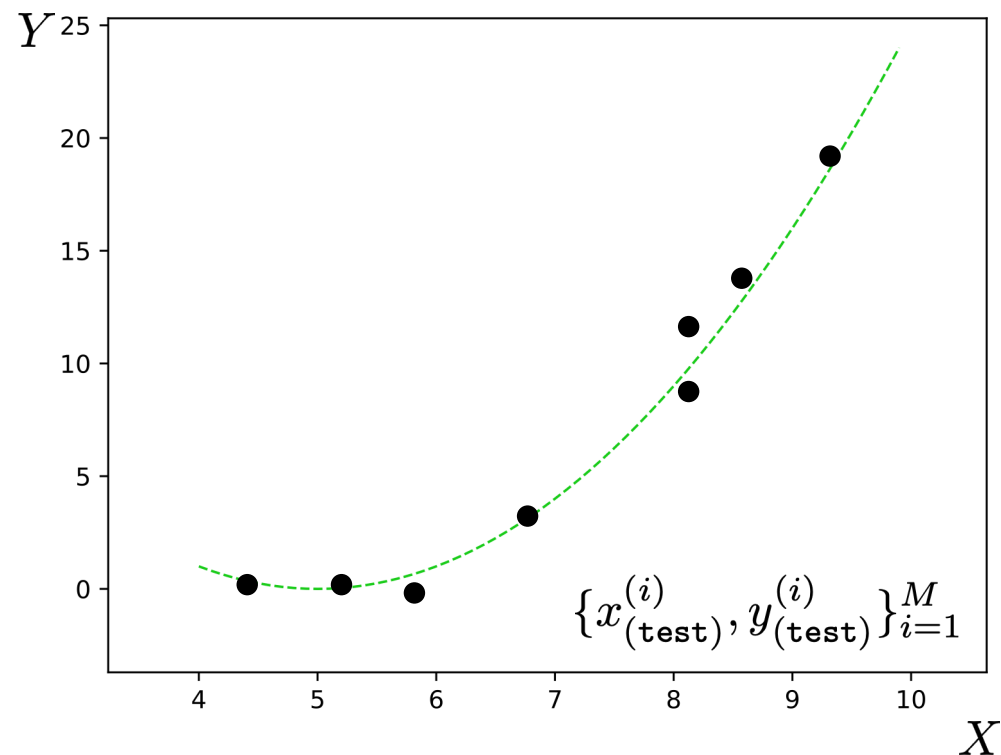
Training data



True data-generating process

p_{data}

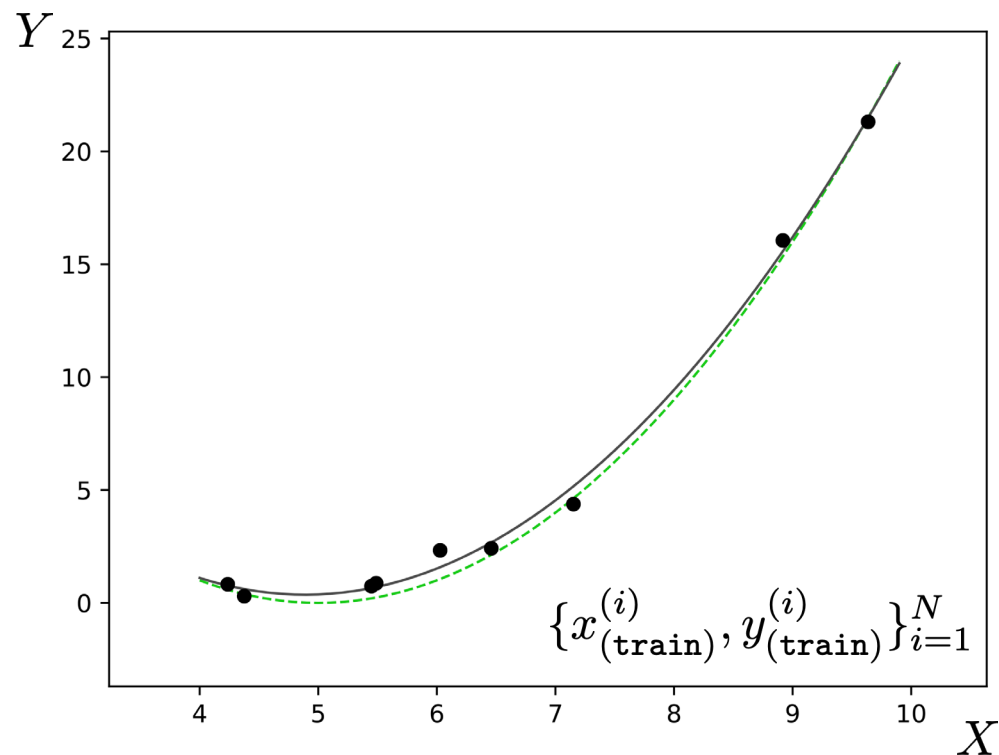
Test data



$$\{x_{(\text{train})}^{(i)}, y_{(\text{train})}^{(i)}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

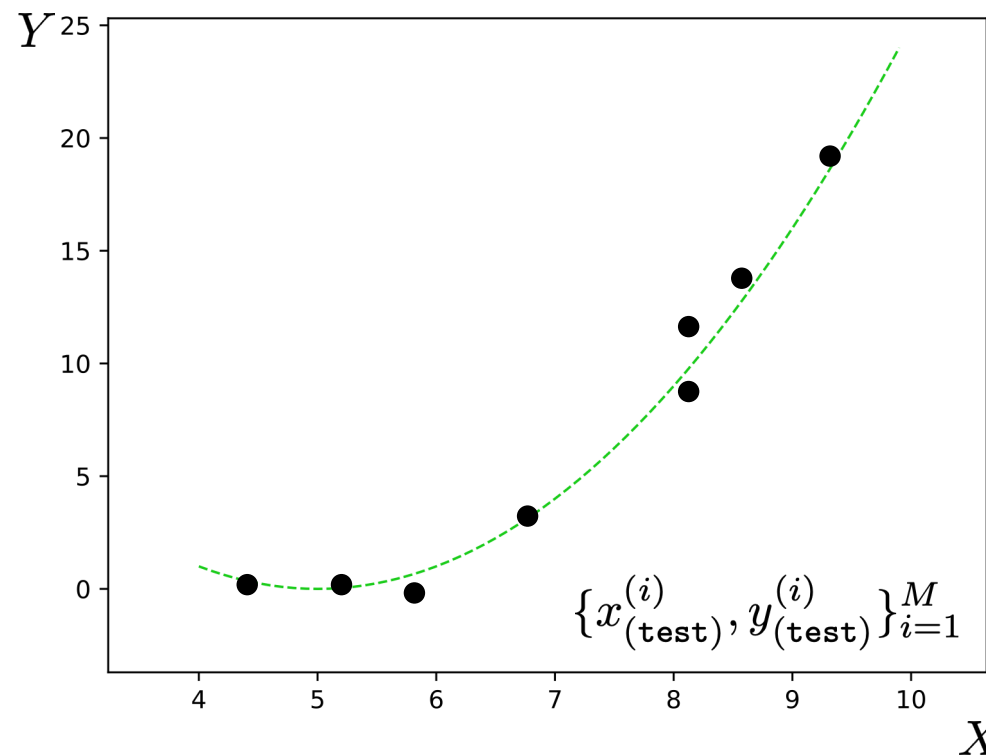
$$\{x_{(\text{test})}^{(i)}, y_{(\text{test})}^{(i)}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

Training data



This is a huge assumption!
Almost never true in practice!

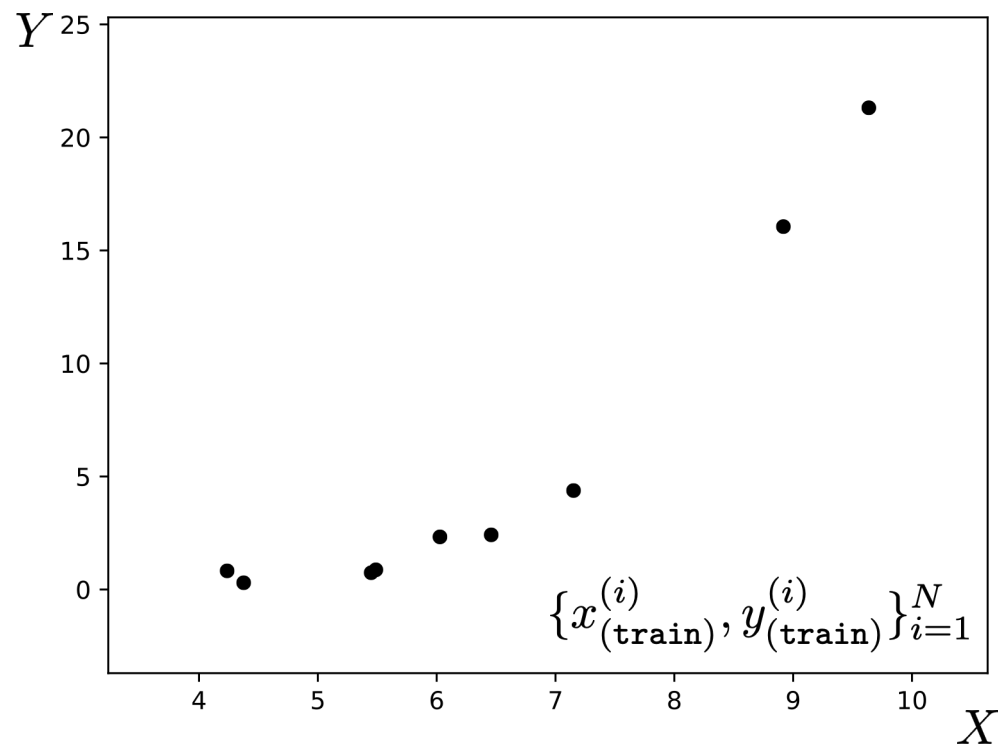
Test data



$$\{x_{(\text{train})}^{(i)}, y_{(\text{train})}^{(i)}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

$$\{x_{(\text{test})}^{(i)}, y_{(\text{test})}^{(i)}\} \stackrel{\text{iid}}{\sim} p_{\text{data}}$$

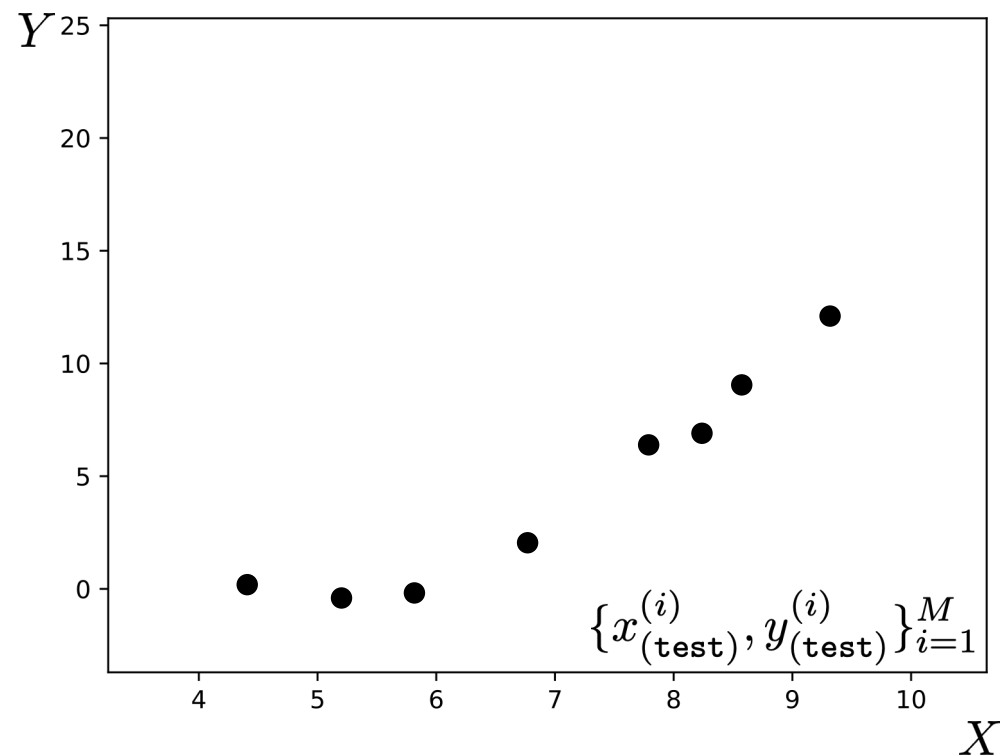
Training data



Much more commonly, we have

$$p_{\text{train}} \neq p_{\text{test}}$$

Test data



$$\begin{aligned} \{x_{(\text{train})}^{(i)}, y_{(\text{train})}^{(i)}\} &\stackrel{\text{iid}}{\sim} p_{\text{train}} \\ \{x_{(\text{test})}^{(i)}, y_{(\text{test})}^{(i)}\} &\stackrel{\text{iid}}{\sim} p_{\text{test}} \end{aligned}$$

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.)

training domain

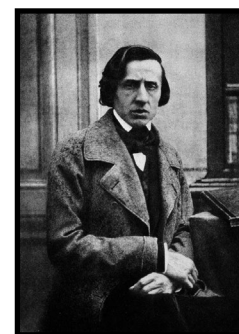
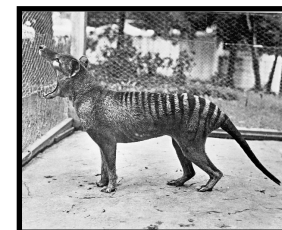
testing domain

(where we actual use our model)

Domain gap between p_{train} and p_{test}
will cause us to fail to generalize.

Space of natural images

Training data



Test data