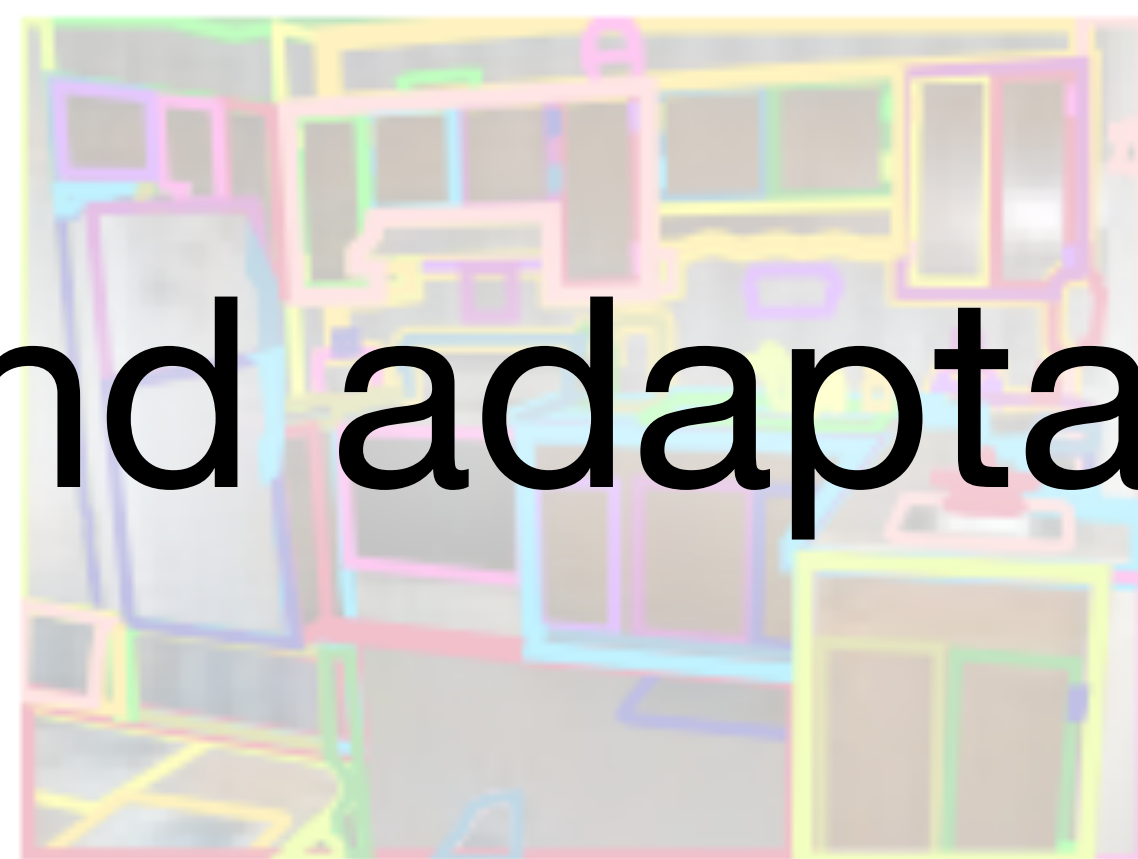# Lecture 24
## Datasets, bias, and adaptation

# Garbage in, garbage out

A machine learning algorithm will do whatever the training data tells it to do.

If the data is bad or biased, the learned algorithm will be too.

# Microsoft's Tay chatbot



Chatbot released on twitter.

Learned from interactions with users (?)

Started mimicking offensive language, was shut down.

what is the yellow thing?

Submit

Predicted top-5 answers with confidence:

frisbee

79.844%

surfboard

7.319%

banana

2.844%

lemon

2.438%

surfboards

1.252%

how many trains are in the picture?

Submit

Predicted top-5 answers with confidence:

3

30.233%

5

18.270%

4

17.000%

2

11.343%

6

7.806%

Of number questions (e.g. "how many…"), 26.04% of the time, the answer is 2
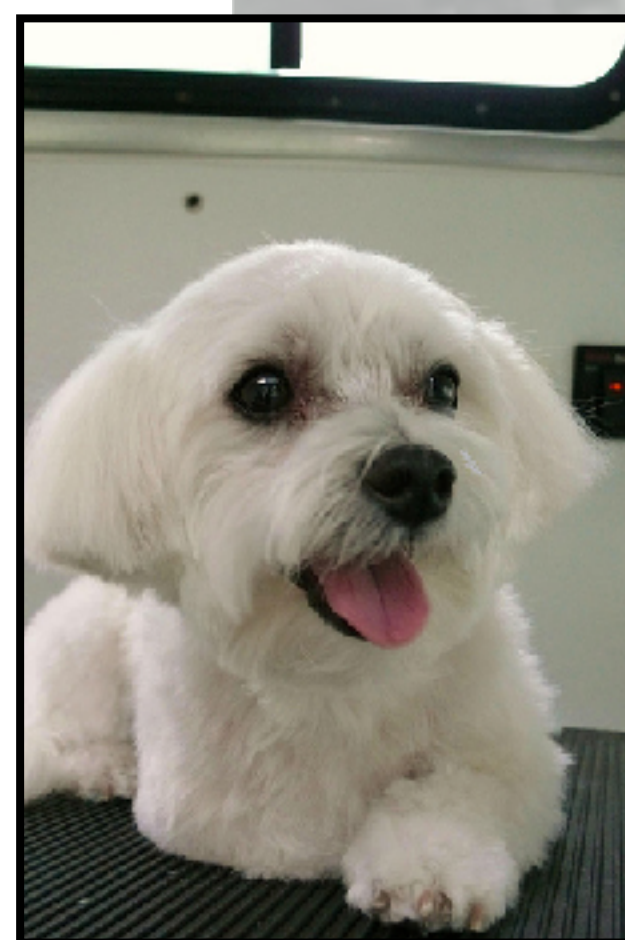
Of yes/no questions, 58.83% of the time, the answer is yes

[VQA, Agrawal, Lu, Antol et al., https://arxiv.org/pdf/1505.00468.pdf

["Colorful image colorization", Zhang et al., ECCV 2016]

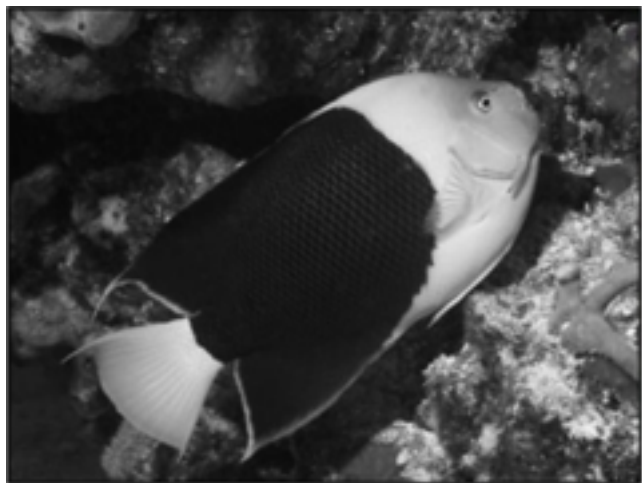["Colorful image colorization", Zhang et al., ECCV 2016]

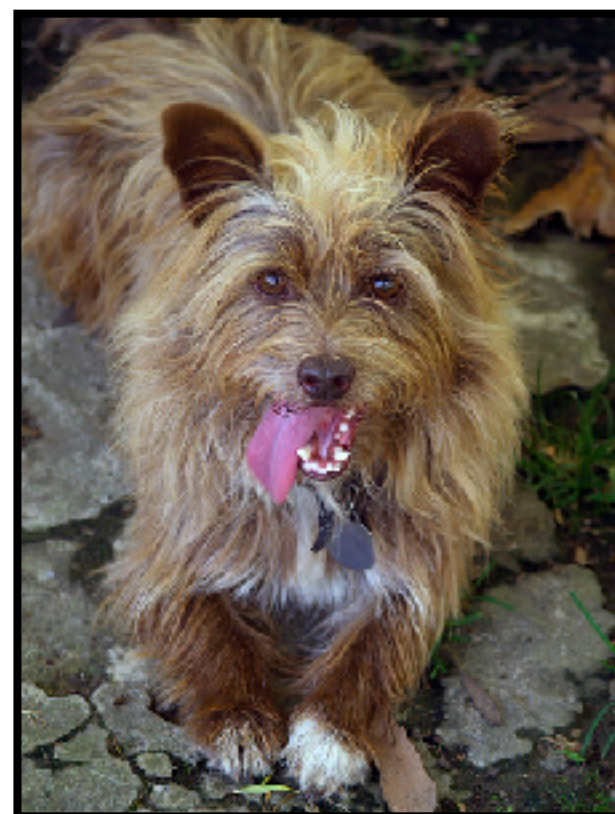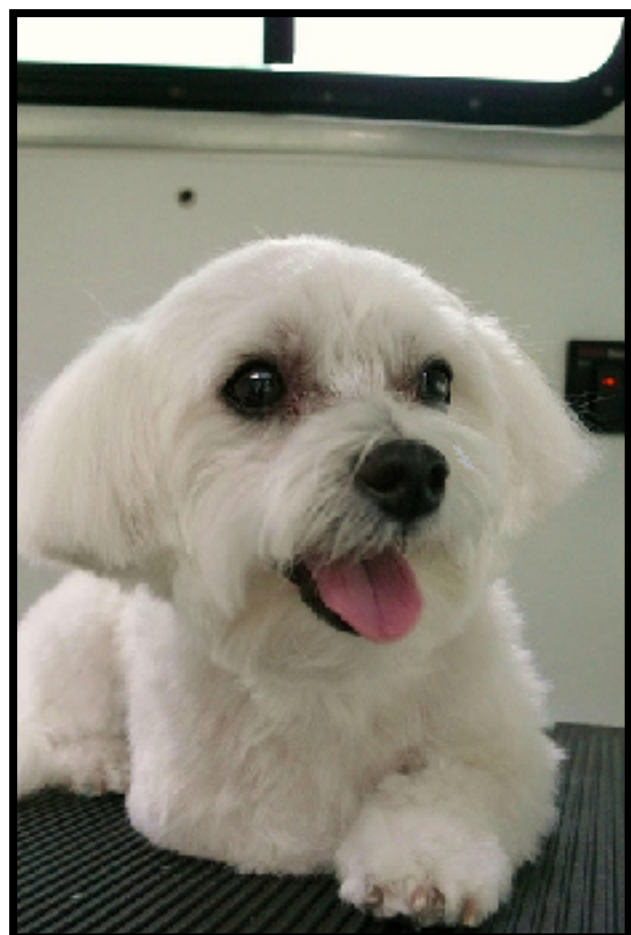["Colorful image colorization", Zhang et al., ECCV 2016]

**x**      **y**

$\mathbf{x}'$

Training data

Test data

# Training data

# Test data

## Driving simulator (GTA)

## Driving in the real world

Let's revisit the problem of generalization

Training data
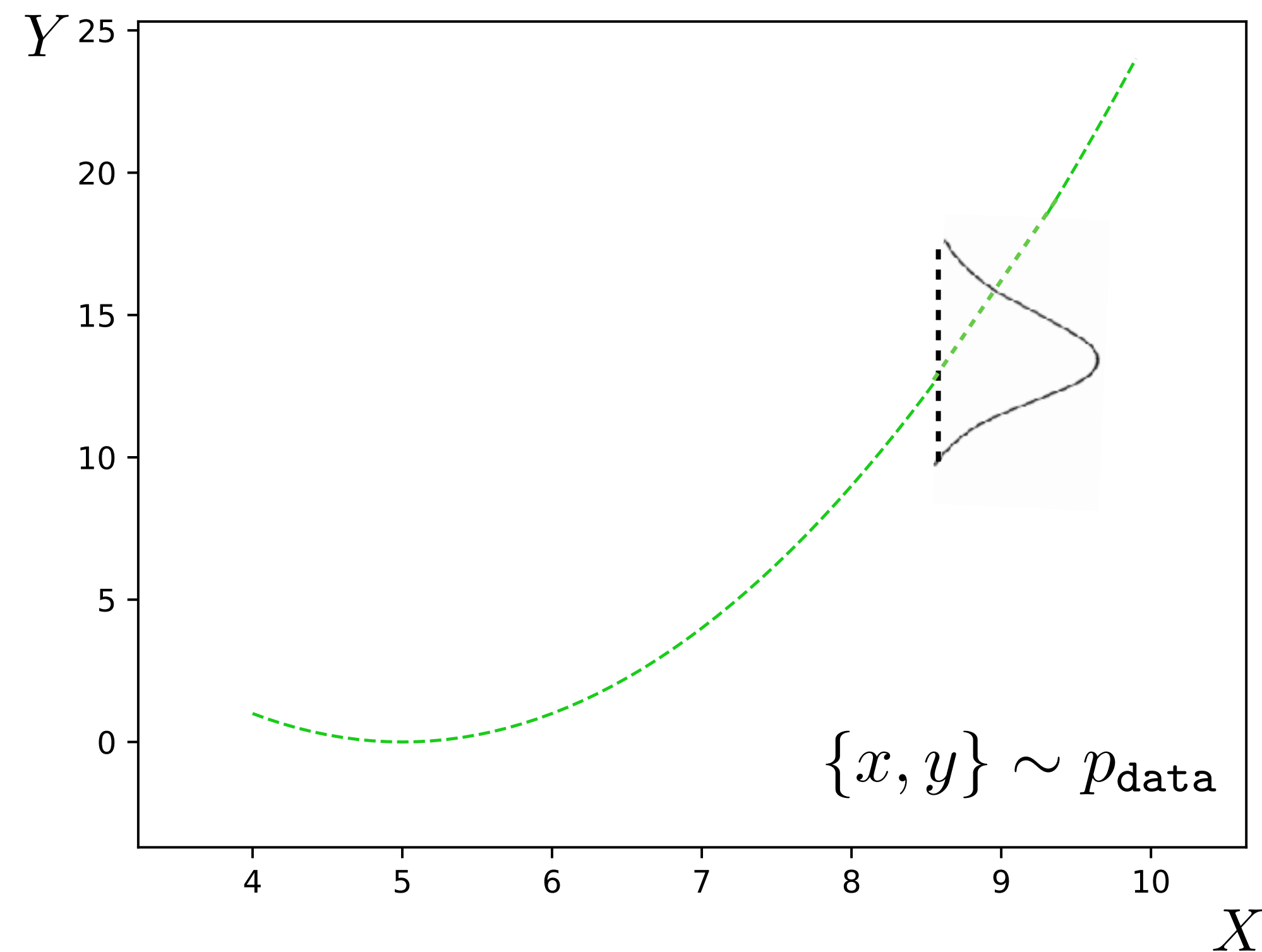
$Y$

$\{x_i, y_i\}_{i=1}^{N}$

$X$

## Training data

$Y$

$\{x_i, y_i\}_{i=1}^N$

$X$

## Test data

$Y$

$\{x, y\} \sim p_{\texttt{data}}$

$X$

True **data-generating process**

$p_{\texttt{data}}$

## Training data

$$Y$$

$$\{x_i^{(\mathtt{train})}, y_i^{(\mathtt{train})}\}_{i=1}^{N}$$

$$X$$

## Test data

$$Y$$

$$\{x_i^{(\mathtt{test})}, y_i^{(\mathtt{test})}\}_{i=1}^{M}$$

$$X$$

True **data-generating process**

$$p_{\mathtt{data}}$$

$$\{x_i^{(\mathtt{train})}, y_i^{(\mathtt{train})}\} \overset{\mathtt{iid}}{\sim} p_{\mathtt{data}}$$

$$\{x_i^{(\mathtt{test})}, y_i^{(\mathtt{test})}\} \overset{\mathtt{iid}}{\sim} p_{\mathtt{data}}$$

Training data

Test data

$\{x_i^{(\text{train})}, y_i^{(\text{train})}\}_{i=1}^N$

$\{x_i^{(\text{test})}, y_i^{(\text{test})}\}_{i=1}^M$

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

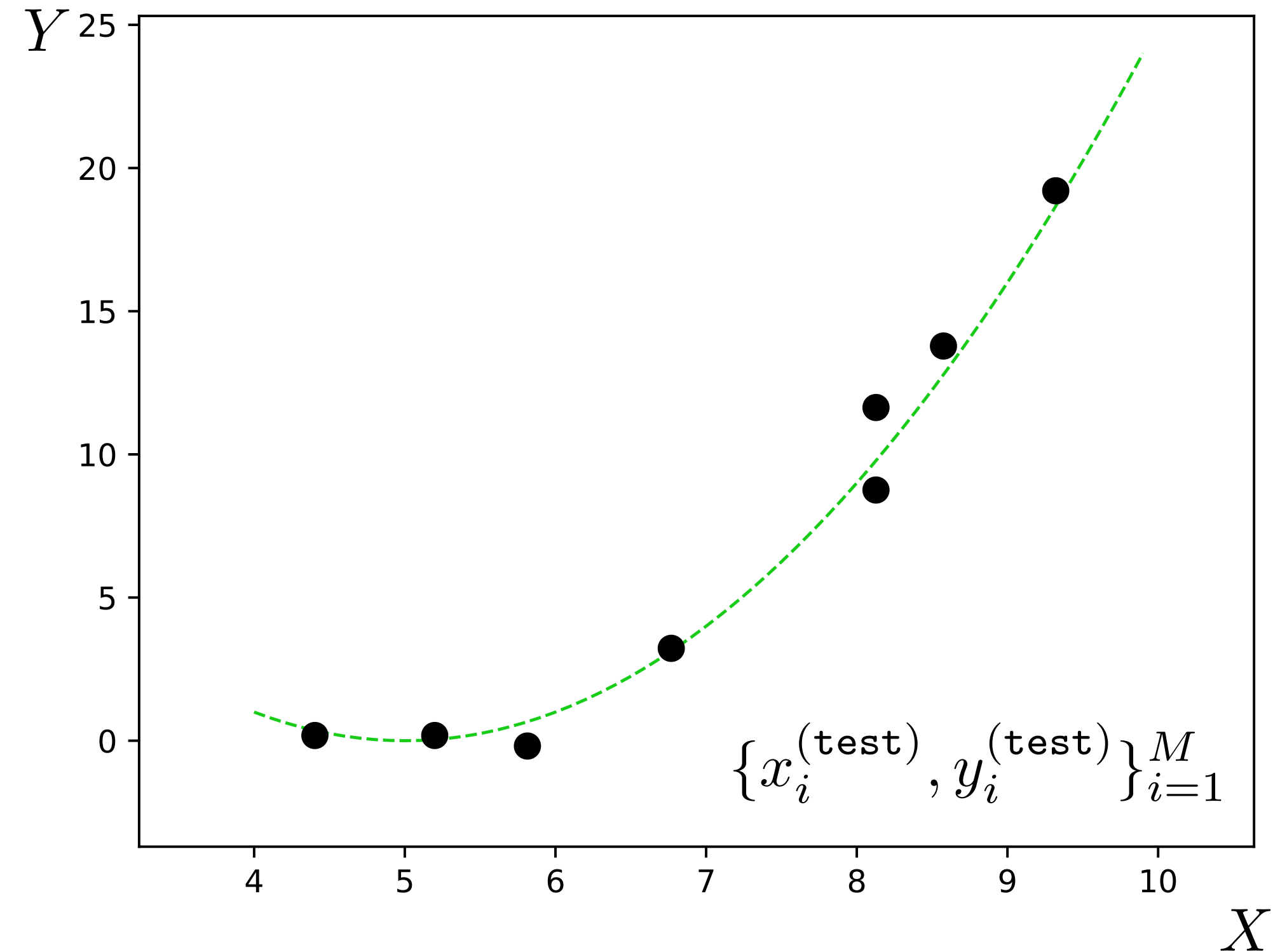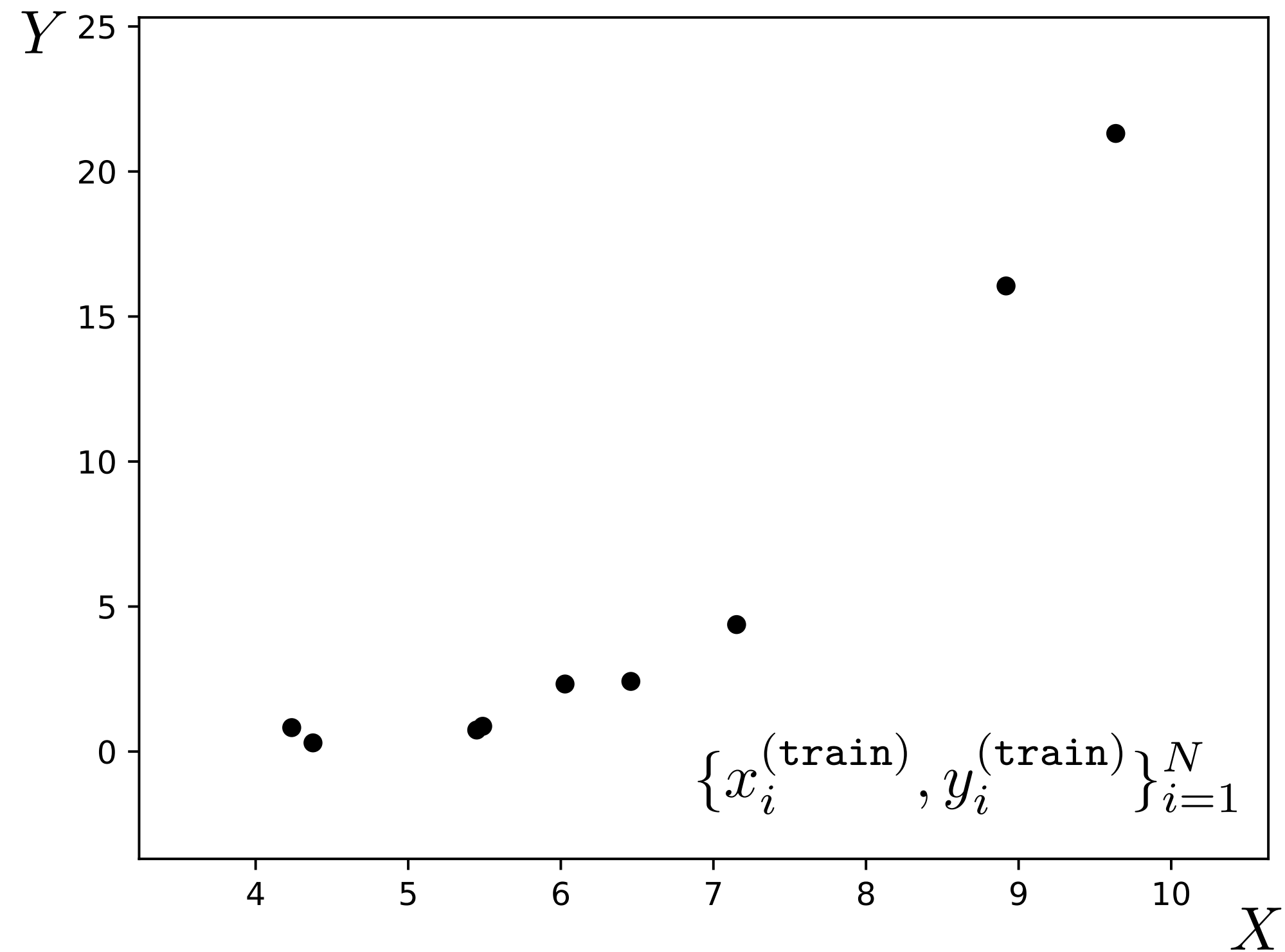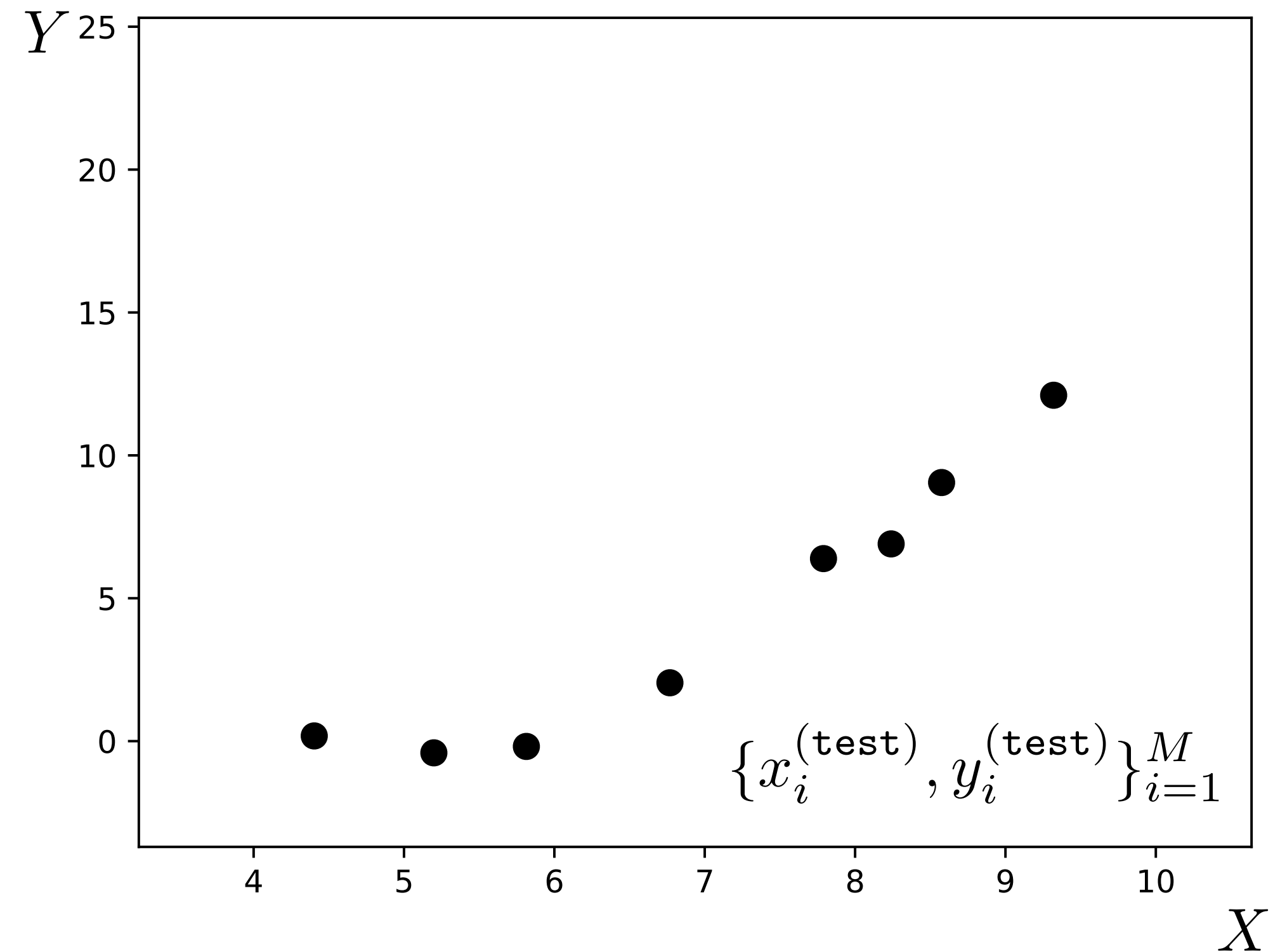$$\{x_i^{(\text{train})}, y_i^{(\text{train})}\} \overset{\text{iid}}{\sim} p_{\text{data}}$$

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

## Training data

$Y$

$\{x_i^{(\texttt{train})}, y_i^{(\texttt{train})}\}_{i=1}^{N}$

$X$

## Test data

$Y$

$\{x_i^{(\texttt{test})}, y_i^{(\texttt{test})}\}_{i=1}^{M}$
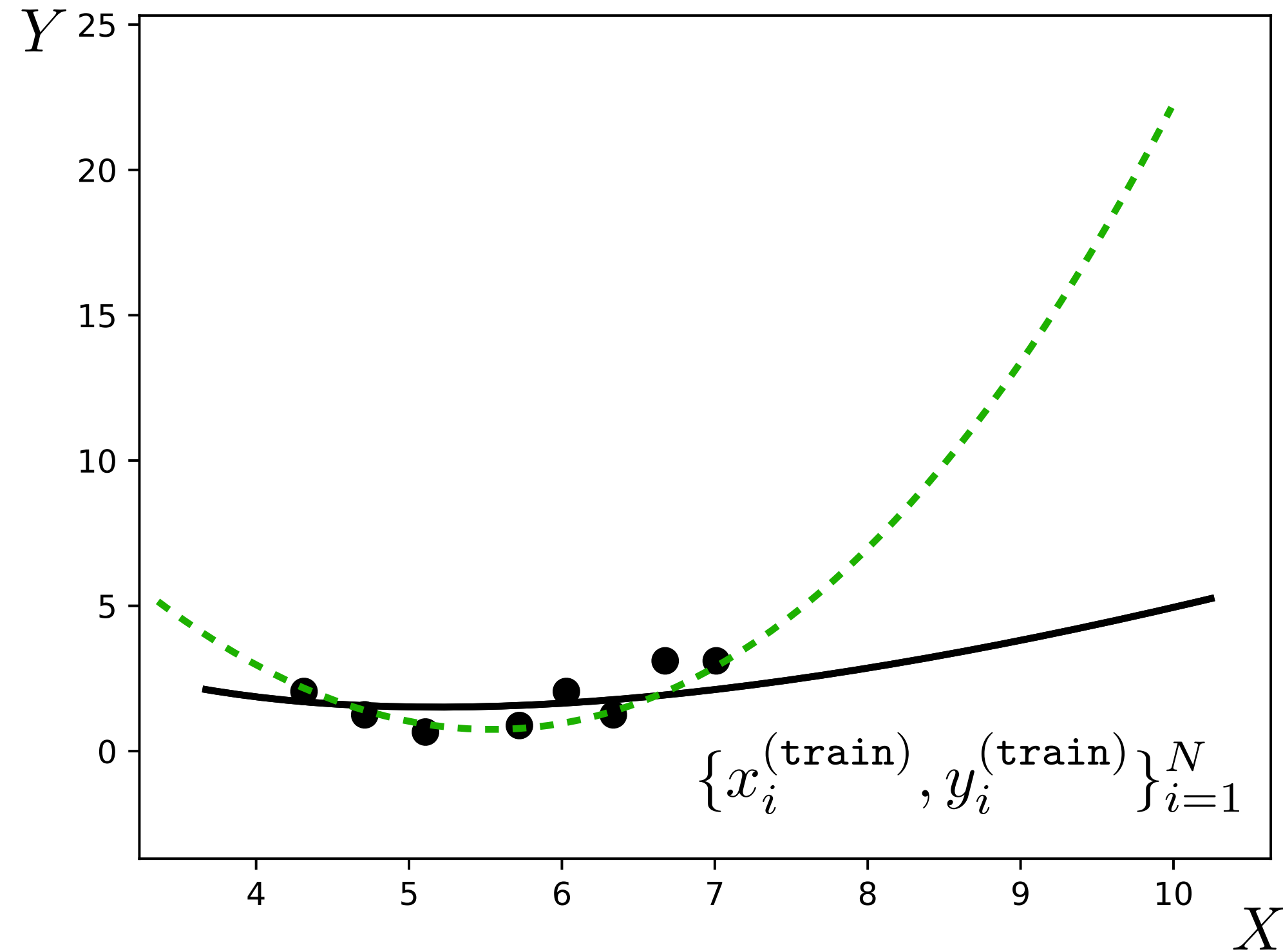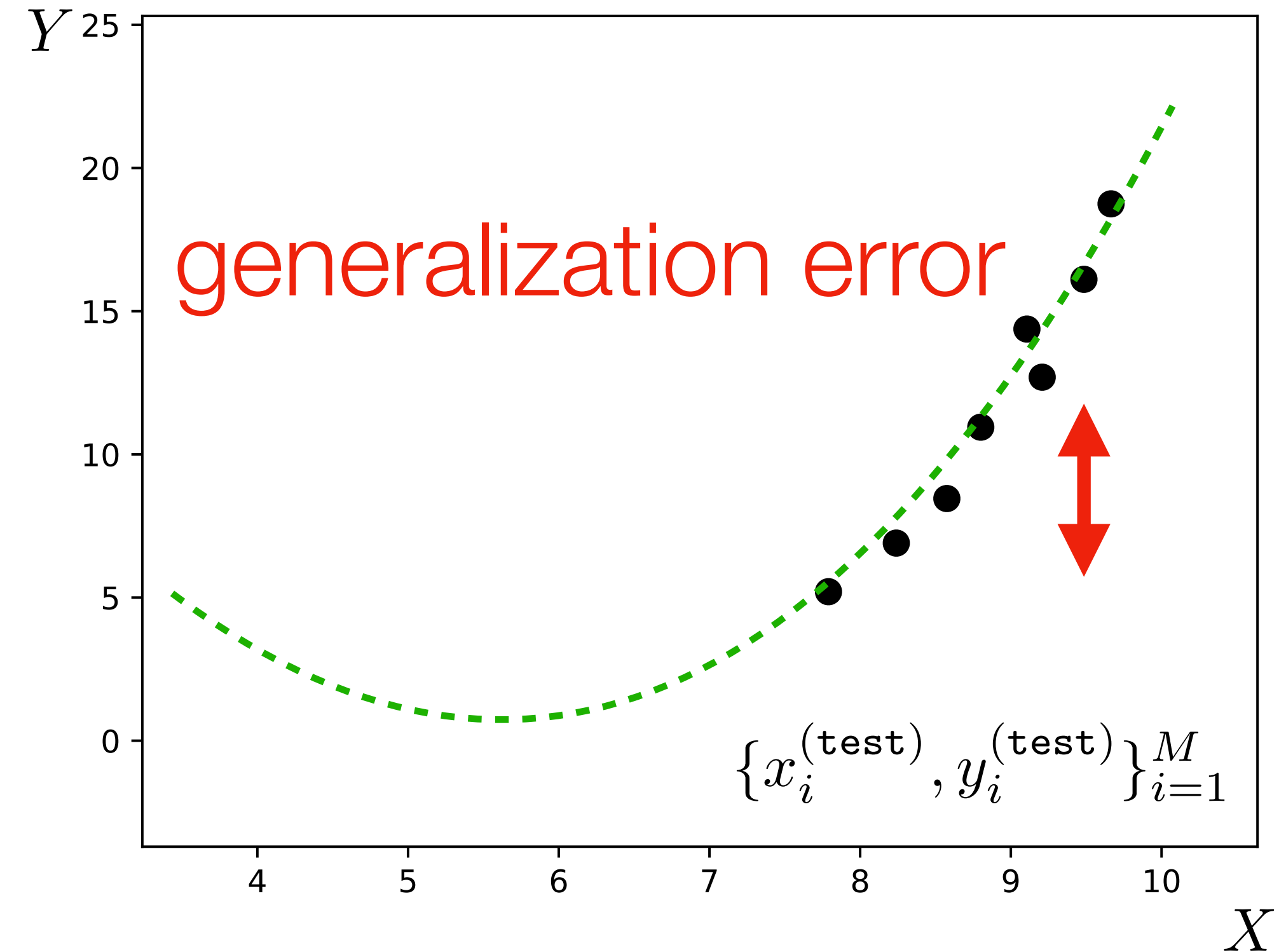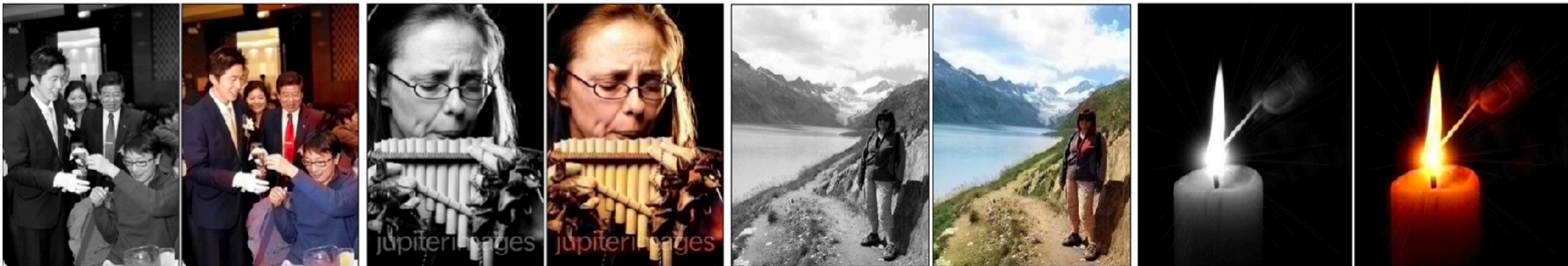
$X$

Much more commonly, we have

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

$$\{x_i^{(\texttt{train})}, y_i^{(\texttt{train})}\} \overset{\texttt{iid}}{\sim} p_{\texttt{train}}$$

$$\{x_i^{(\texttt{test})}, y_i^{(\texttt{test})}\} \overset{\texttt{iid}}{\sim} p_{\texttt{test}}$$

Training data

Test data

$Y$

$\{x_i^{(\text{train})}, y_i^{(\text{train})}\}_{i=1}^{N}$

$\{x_i^{(\text{test})}, y_i^{(\text{test})}\}_{i=1}^{M}$

generalization error
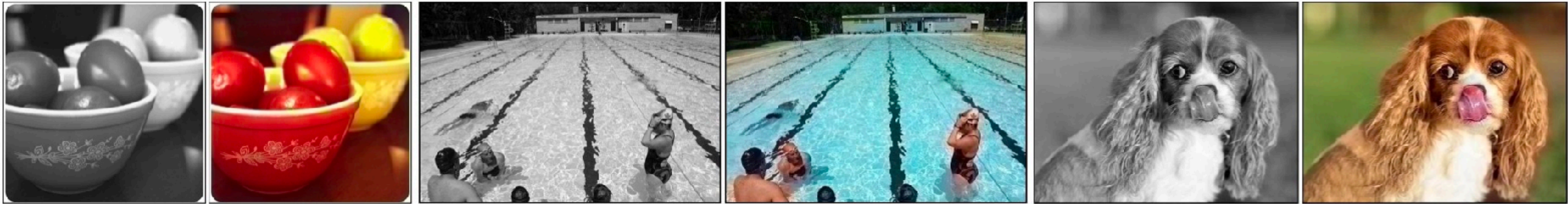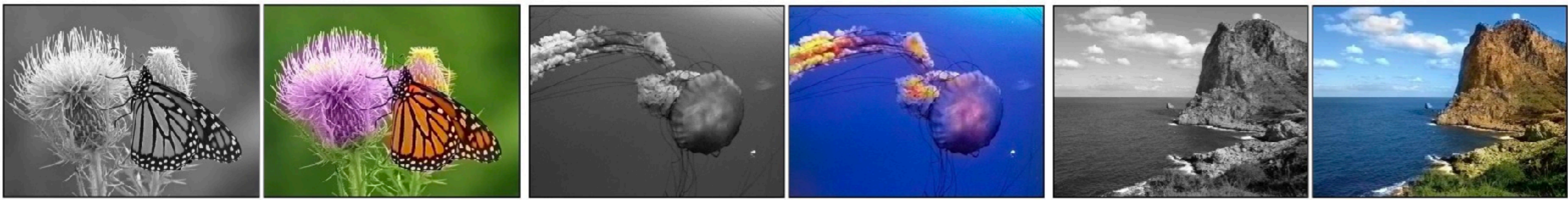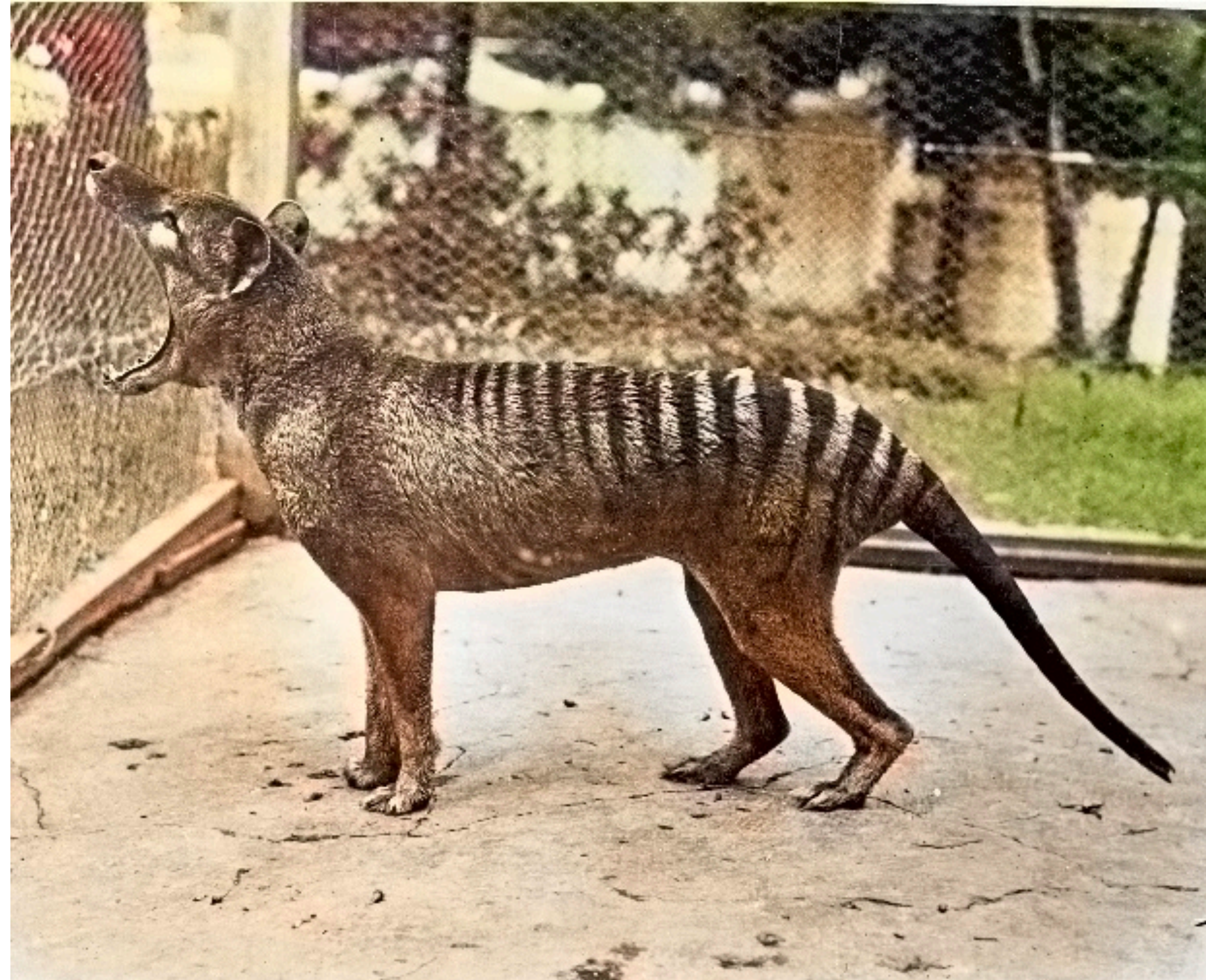
$X$

Our training data did cover the part of the distribution that was tested
(**biased data**)

u/Rafael_P_S



Thylacine



Chopin

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



Space of natural images

Training data

Test data

# Algorithmic Bias



http://gendershades.org/overview.html

http://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf

While this study focused on gender classification, the machine learning techniques used to determine gender are also broadly applied to many other areas of facial analysis and automation. Face recognition technology that has not been publicly tested for demographic accuracy is increasingly used by law enforcement and at airports. AI fueled automation now helps determine who is fired, hired, promoted, granted a loan or insurance, and even how long someone spends in prison.

For interested readers, authors Cathy O'Neil and Virginia Eubanks explore the real-world impact of algorithmic bias.



"This book is downright scary—but...you will emerge smarter and more empowered to demand justice." —NAOMI KLEIN

**AUTOMATING INEQUALITY**

HOW HIGH-TECH TOOLS PROFILE, POLICE, AND PUNISH THE POOR

VIRGINIA EUBANKS



**WEAPONS OF MATH DESTRUCTION**

HOW BIG DATA INCREASES INEQUALITY AND THREATENS DEMOCRACY

CATHY O'NEIL

http://gendershades.org/overview.html

# How can we collect good data?

+ Correctly labeled
+ Unbiased (good coverage of all relevant kinds of data)

Crowdsourcing

# The value of data



The Large Hadron Collider

$ 10^{10}$



Amazon Mechanical Turk

$ 10^{2} - 10^{4}$

# But can humans collect good data?

# Getting more humans in the annotation loop

Labeling to get a Ph.D.



## Labeling for fun
Luis Von Ahn and Laura Dabbish 2004



## Labeling for money
(Sorokin, Forsyth, 2008)


amazonmechanical turk
Artificial Artificial Intelligence
beta

Labeling because it
gives you added value



Visipedia
(Belongie, Perona, et al)

Just for labeling


LabelMe

# Beware of the human in your loop

- What do you know about them?

- Will they do the work you pay for?

Let's check a few simple experiments

# People have biases…

Turkers were offered 1 cent to pick a number from 1 to 10.

~850 turkers



Experiment by Greg Little
From http://groups.csail.mit.edu/uid/deneme/

# Do humans have consistent biases?

Choose Item
Requester: SimpleSphere     Reward: $0.01 per HIT     HITs Available: 1     Duration: 60 minutes
Qualifications Required: None

Please choose one of the following:

Results form 100 HITS:



Experiment by Greg Little
From http://groups.csail.mit.edu/uid/deneme/

# Do humans do what you ask for?

**Flip a coin**
 Requester: ROBERT C MILLER          Reward: $0.01 per HIT      HITs Available: 3      Duration: 5 minutes
 Qualifications Required:  None

**Please flip an actual coin and type either H or T below.**

After 50 HITS:

And 50 more:

31 heads, 19 tails

34 heads, 16 tails

Experiment by Rob Miller

From http://groups.csail.mit.edu/uid/deneme/

# Are humans reliable even in simple tasks?

Choose the given item.
Requester: SimpleSphere        Reward: $0.01 per HIT        HITs Available: 1        Duration: 60 minutes
Qualifications Required: None

Please click button B:

B

C

A

Results of 100 HITS:
A: 2
B: 96
C: 2

Experiment by Greg Little
From http://groups.csail.mit.edu/uid/deneme/

So we can sometimes collect good training data.

But suppose we messed up. Our test setting doesn't look like the training data!

How can we bridge the domain gap?

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



Space of natural images

Training data

Test data

*source domain*

*target domain*
(where we actual use our model)

**Domain gap** between $p_{\texttt{source}}$ and $p_{\texttt{target}}$ will cause us to fail to generalize.

Space of natural images

Source data

Target data

# Idea #1: transform the target domain to look like the source domain

Data space

source data

target data

(Or vice versa)     This is called **domain adaptation**

# Domain adaptation

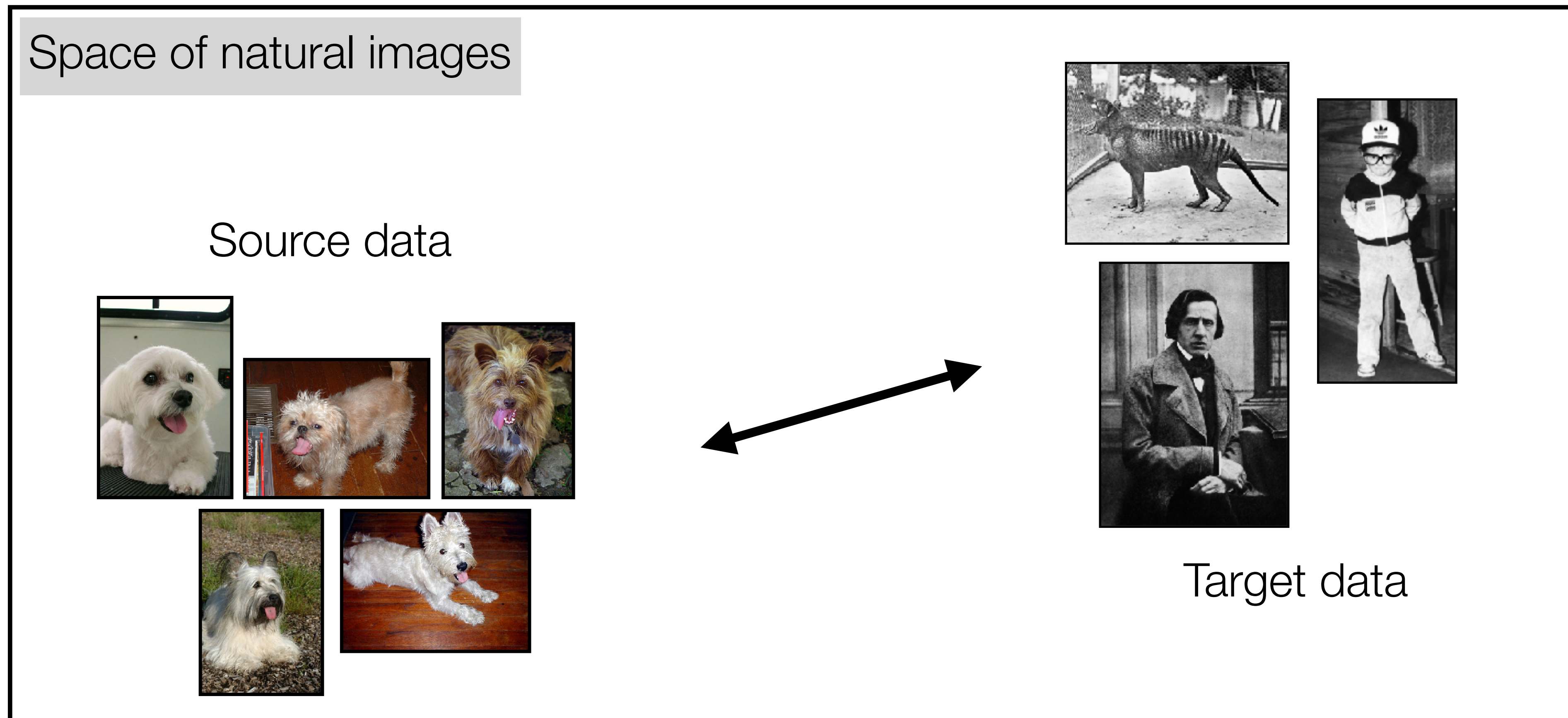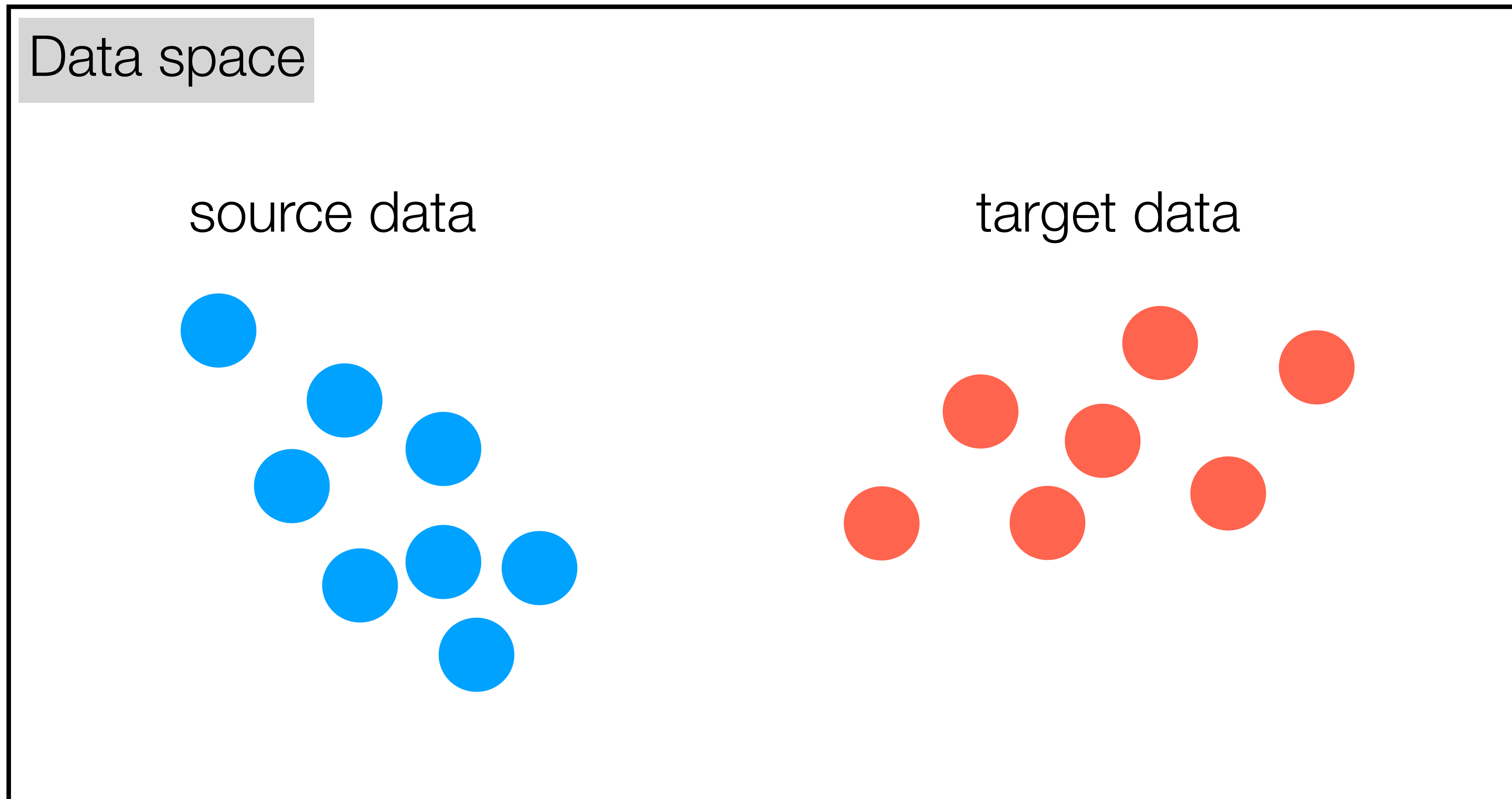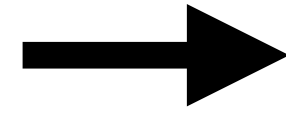- We have source domain pairs $\{\mathbf{x}^{\mathbf{source}}, \mathbf{y}^{\mathbf{source}}\}$

- Learn a mapping F: $\mathbf{x}^{\mathbf{source}} \longrightarrow \mathbf{y}^{\mathbf{source}}$

- We want to apply F to target domain data $\mathbf{x}^{\mathbf{target}}$

- Find transformation T: $\mathbf{x}^{\mathbf{target}} \longrightarrow \mathbf{x}^{\mathbf{source}}$

- Now apply F(T($\mathbf{x}^{\mathbf{target}}$)) to predict $\mathbf{y}^{\mathbf{target}}$

$p_{\texttt{source}}$

$p_{\texttt{target}}$

It's a just another distribution mapping problem!

# GANs

Gaussian

Target distribution

**z**

**Y**

# CycleGAN

Horses

Zebras



X

→

Y

# Domain adaptation



$p_{\texttt{source}}$ $\rightarrow$ $p_{\texttt{target}}$

**Domain gap** between $p_{\texttt{source}}$ and $p_{\texttt{target}}$ will cause us to fail to generalize.
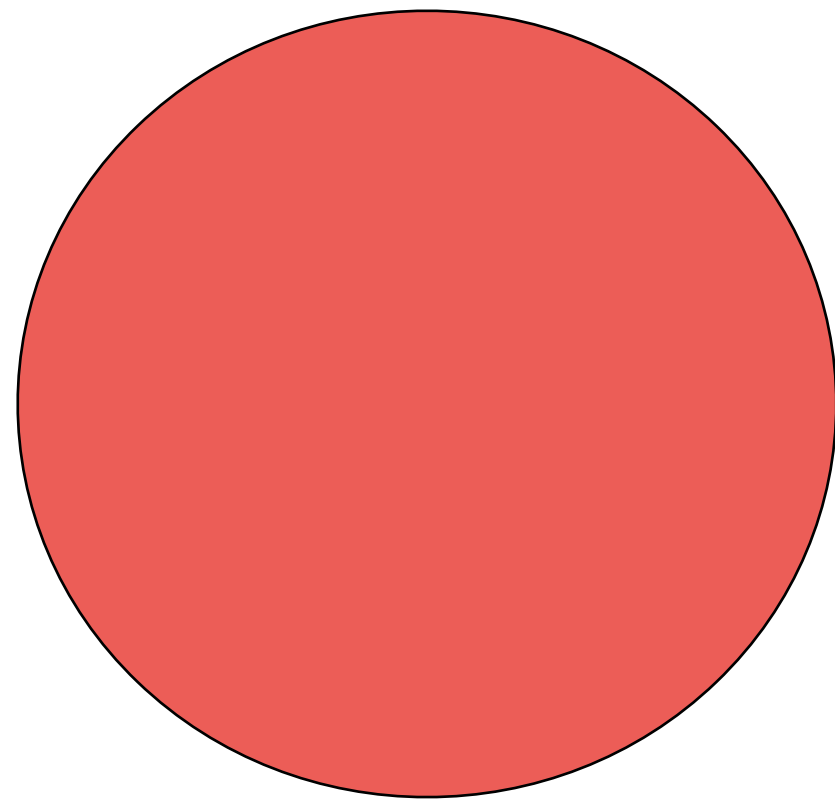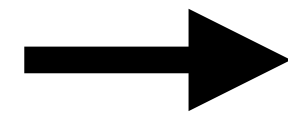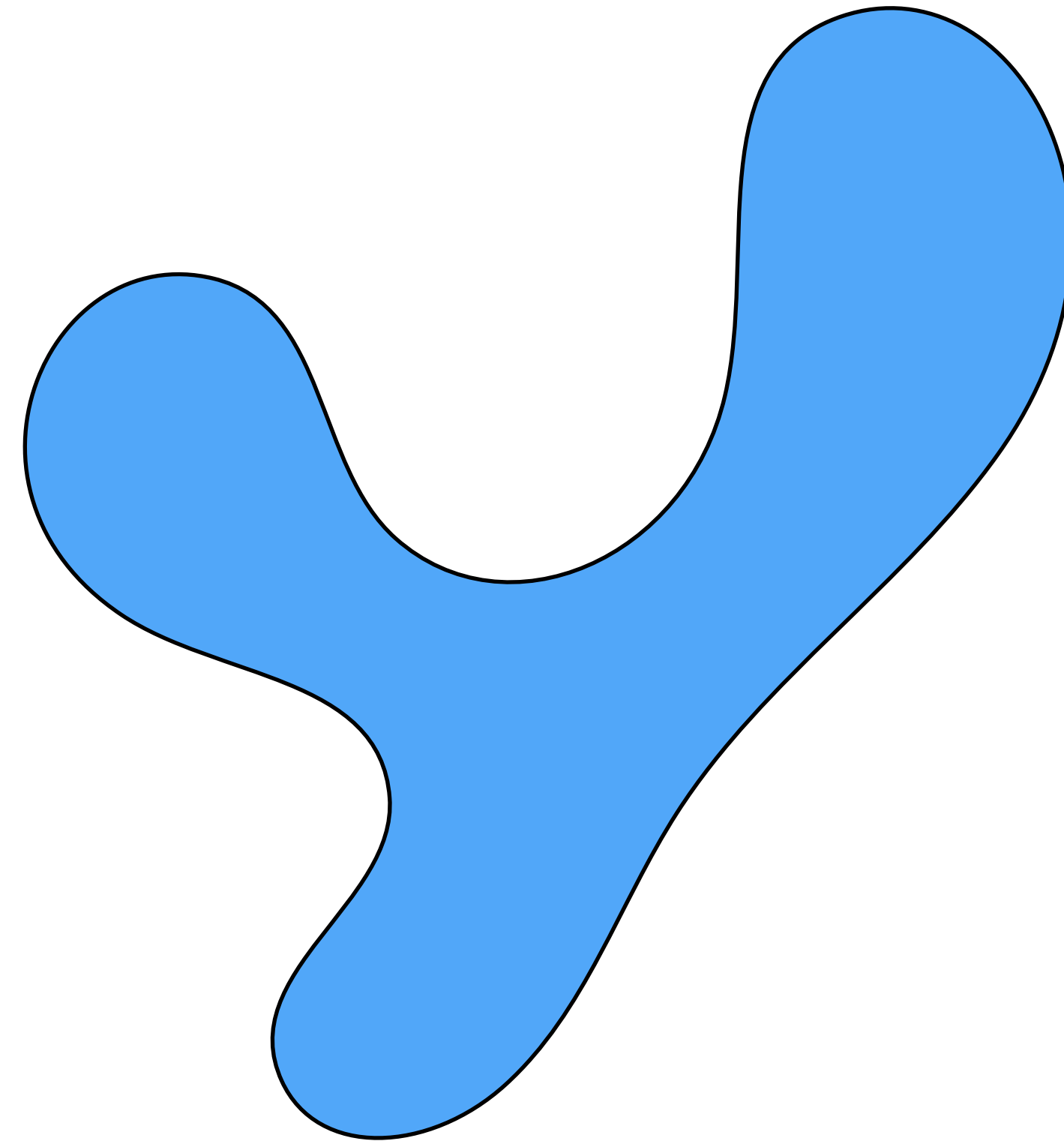


Space of images

Source data

Target data

# CyCADA: Cycle-Consistent Adversarial Domain Adaptation

Source domain

Target domain



[Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Darrell, Efros, arXiv 2017]

# CycleGAN



Training data

# CycleGAN



Training data

CycleGAN                    FCN



Training data

# OpenAI Dactyl



FINGER PIVOTING

SLIDING

FINGER GAITING

Idea #2: train on randomly perturbed data, so that test set just looks like another random perturbation

Data space

Training data

Test data

This is called **domain randomization** or **data augmentation**

# Data augmentation



*Training data*

$\mathbf{x}$    $y$

"Fish" ,

"Grizzly" ,

"Chameleon" ,

⋮

"Fish"    Mirror

"Fish"    Crop

"Fish"    Crop

"Fish"    Darken

# Domain randomization

Training data

Test data



[Sadeghi & Levine 2016]

Above example is from [Tobin et al. 2017]

Table 1: Ranges of physics parameter randomizations.

| Parameter | Scaling factor range | Additive term range |
|---|---|---|
| object dimensions | uniform($[0.95, 1.05]$) | |
| object and robot link masses | uniform($[0.5, 1.5]$) | |
| surface friction coefficients | uniform($[0.7, 1.3]$) | |
| robot joint damping coefficients | loguniform($[0.3, 3.0]$) | |
| actuator force gains (P term) | loguniform($[0.75, 1.5]$) | |
| joint limits | | $\mathcal{N}(0, 0.15)$ rad |
| gravity vector (each coordinate) | | $\mathcal{N}(0, 0.4)$ m/s$^2$ |





- All Randomizations
- No Randomizations

What if we go waaaay outside of the training distribution?

Our training data did not cover the part of the distribution that was tested (**biased data**)

Data space

Training data

Test data

*Out here, model response is highly unpredictable*

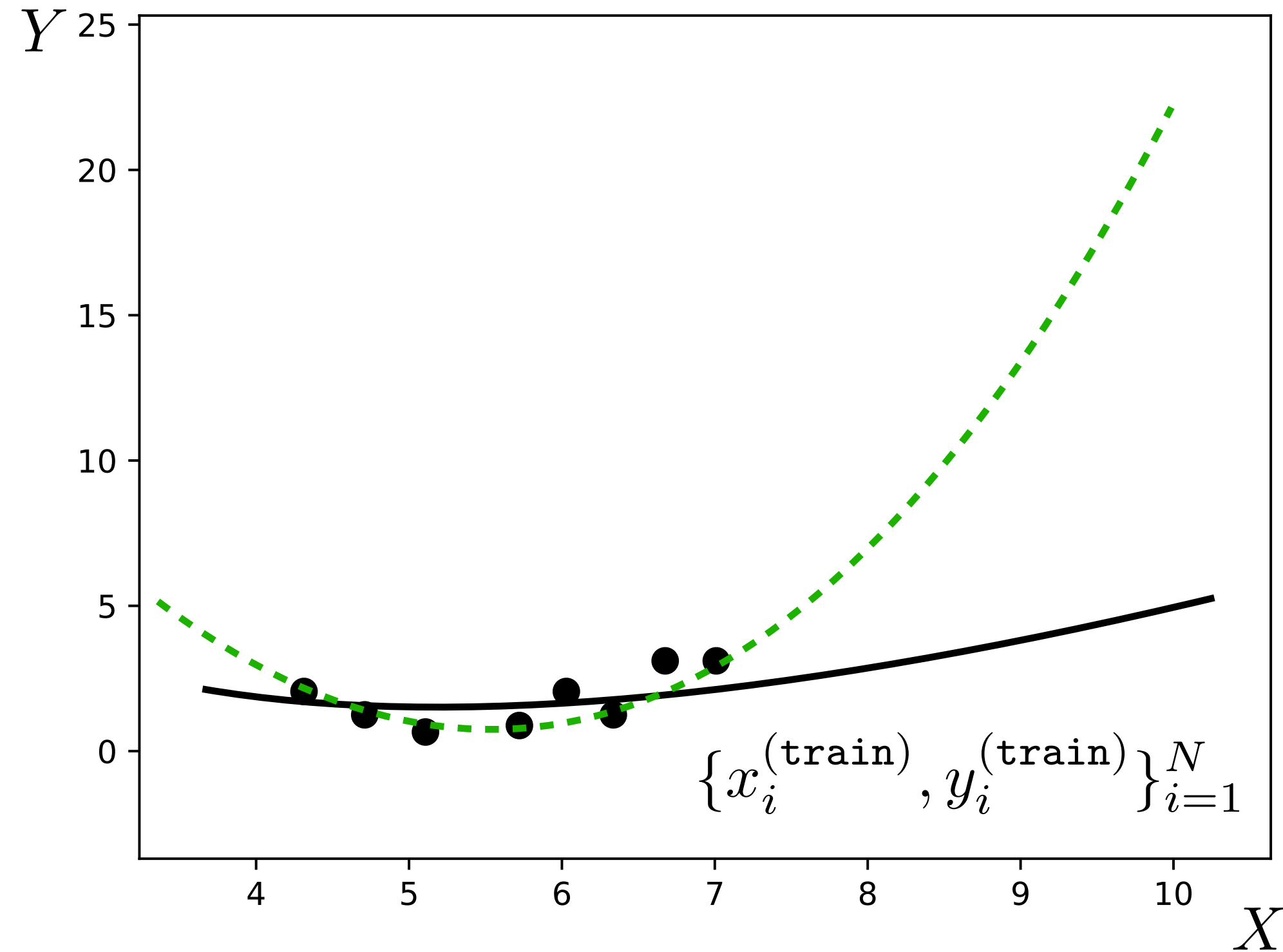# "Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images"
[Nguyen, Yosinski, and Clune, CVPR 2015]

# "Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images"
## [Nguyen, Yosinski, and Clune, CVPR 2015]

# Weirdness of high-dimensional space:

Data space

Training data



Usually, there are *blind spots* where the model has not fit the distribution well, and behaves unpredictably
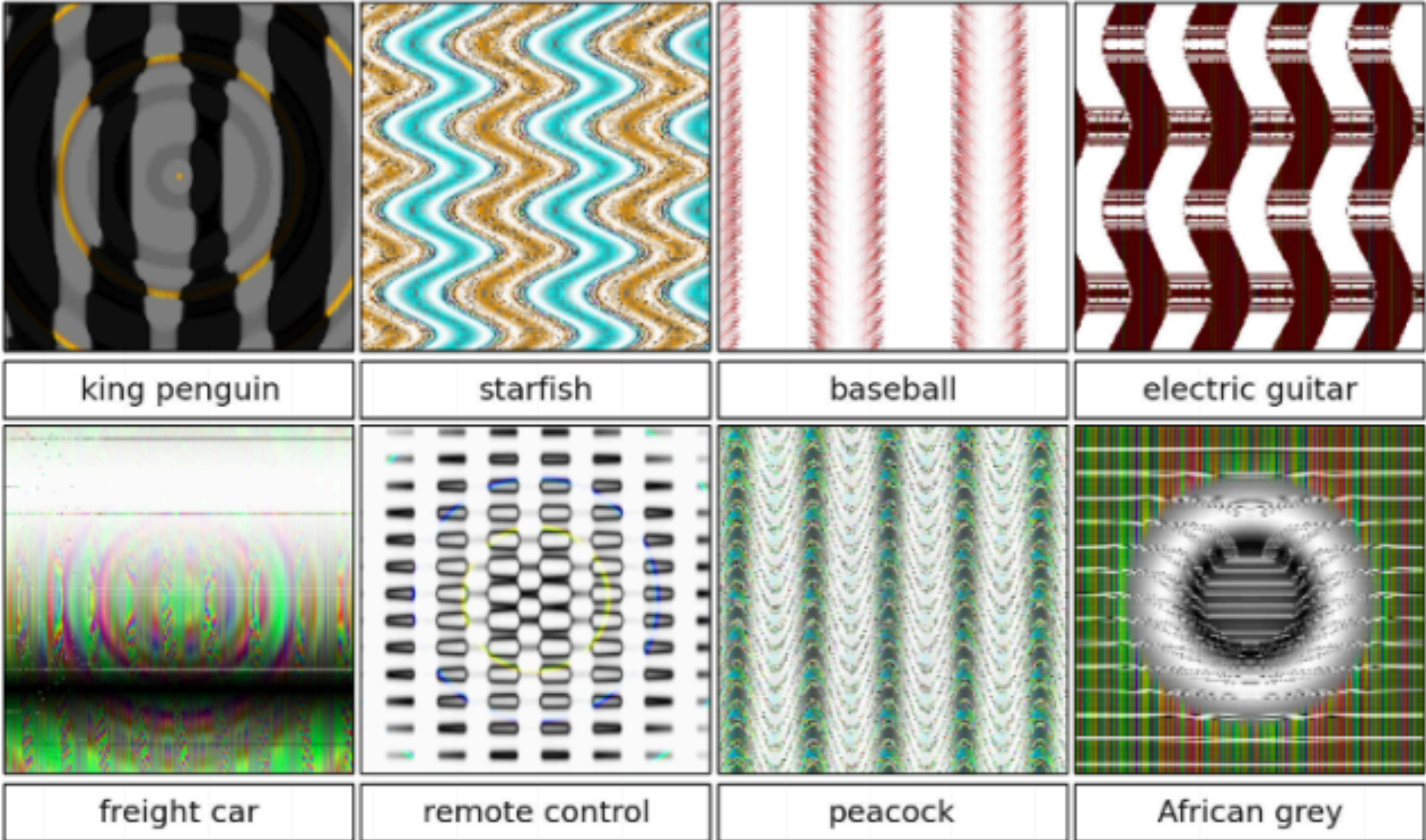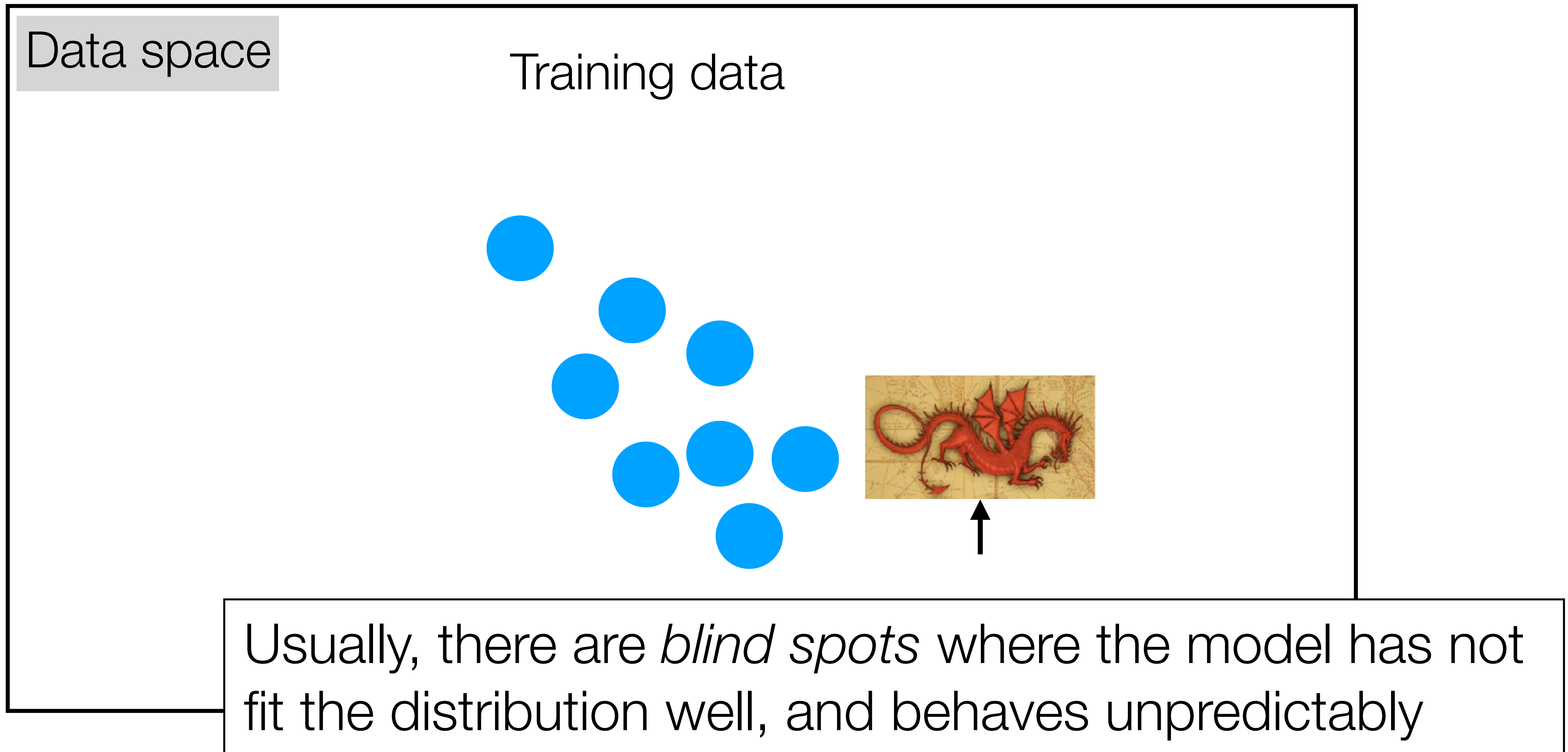
# Adversarial noise

**x**



**+**

**r**



**=**

**x + r**



$f$

$f$

$y$       "School bus"

"Ostrich"

$$\arg \max_{\mathbf{r}} p(y = \texttt{ostrich}|\mathbf{x} + \mathbf{r}) \quad \text{subject to} \quad \|\mathbf{r}\| < \epsilon$$

["Intriguing properties of neural networks", Szegedy et al. 2014]

# Anything to worry about?

"NO Need to Worry about Adversarial Examples in Object Detection in Autonomous Vehicles", Lu et al. 2017



(Early) 2017's attacks fail on physical objects, since they are optimized to attack a single view!

# Anything to worry about?

Later in 2017…

"Synthesizing Robust Adversarial Examples", Athalye, Engstrom, Ilyas, Kwok, 2017

3D-printed **turtle** model classified as **rifle** from most viewpoints
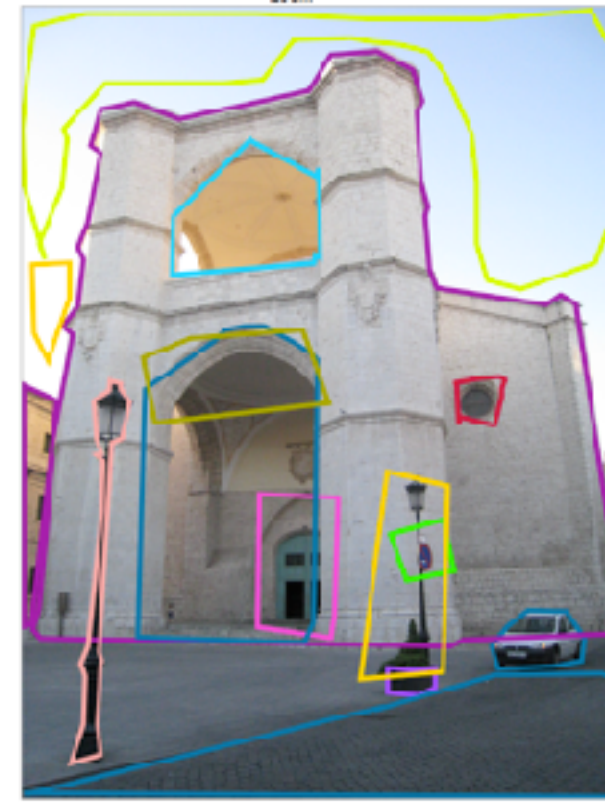
# Anything to worry about?

- Current deep models have bad **worst-case performance**

- Can be exploited by an adversary

- Few guarantees, can't fully trust what the model's output

# Anything else to worry about?

- Our datasets are often poorly labeled

- And usually biased (overrepresent certain categories)

- ML method perform beautifully on laboratory data, but often generalize poorly to real-world data

- Can have negative social consequences