

# 6.036: Introduction to Machine Learning

**Lecture start:** Tuesdays 9:35am

**Lecture team:** Prof. Tamara Broderick & TA Elizabeth Zou

**Questions?** Ask on Piazza

**Materials:** slides, video will all be available on Canvas

**Live Zoom feed:** <https://mit.zoom.us/j/94238622313>

## Today's Plan

- I. Meet the team
- II. Machine learning setup
- III. Linear regression
- IV. Regularization

# 6.036: Introduction to Machine Learning



**Lecture start:** Tuesdays 9:35am

**Lecture team:** Prof. Tamara Broderick & TA Elizabeth Zou

**Questions?** Ask on Piazza

**Materials:** slides, video will all be available on Canvas

**Live Zoom feed:** <https://mit.zoom.us/j/94238622313>

## Today's Plan

- I. Meet the team
- II. Machine learning setup
- III. Linear regression
- IV. Regularization

# 6.036: Introduction to Machine Learning



**Lecture start:** Tuesdays 9:35am

**Lecture team:** Prof. Tamara Broderick & TA Elizabeth Zou

**Questions?** Ask on Piazza: “lecture (week) 2” folder

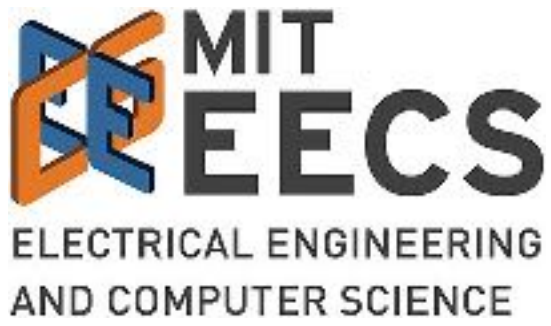
**Materials:** slides, video will all be available on Canvas

**Live Zoom feed:** <https://mit.zoom.us/j/94238622313>

## Today's Plan

- I. Meet the team
- II. Machine learning setup
- III. Linear regression
- IV. Regularization



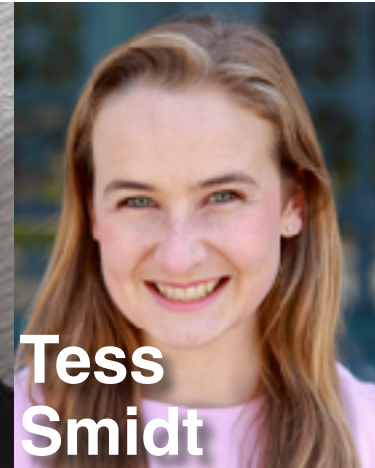


# 6.036: Introduction to Machine Learning, Staff



# 6.036: Introduction to Machine Learning, Staff

Instructors:

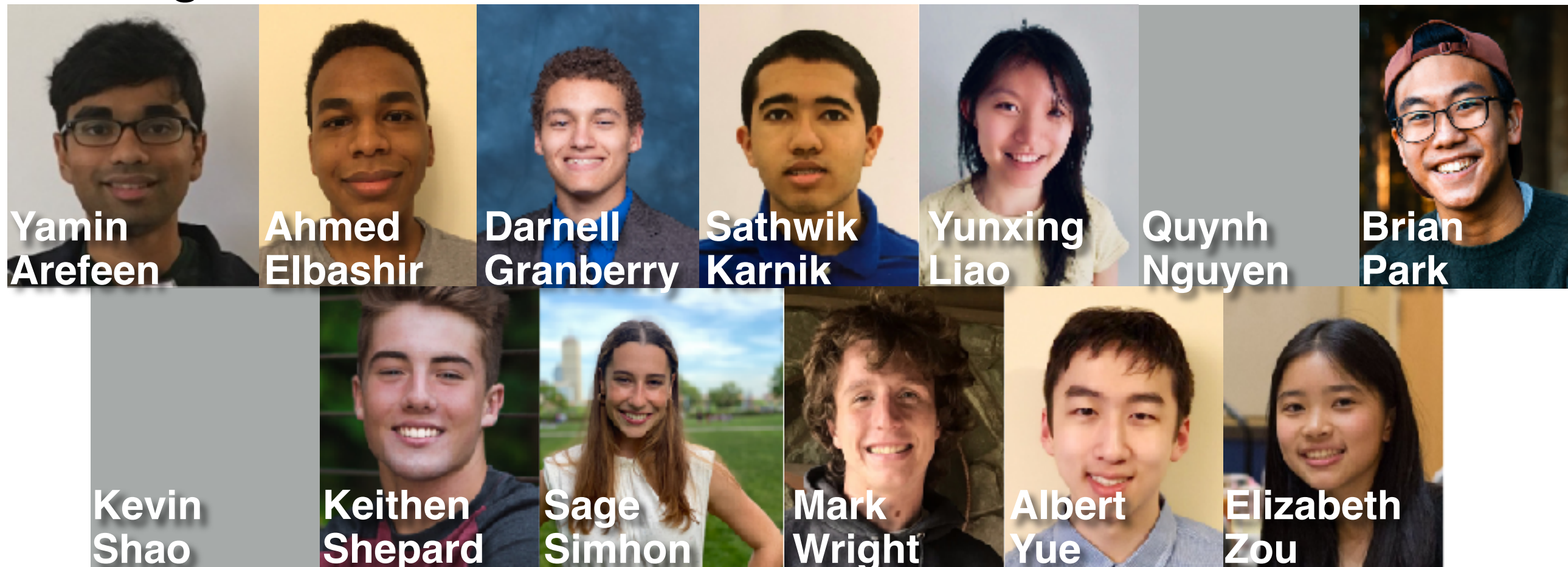


# 6.036: Introduction to Machine Learning, Staff

Instructors:



Teaching Assistants:



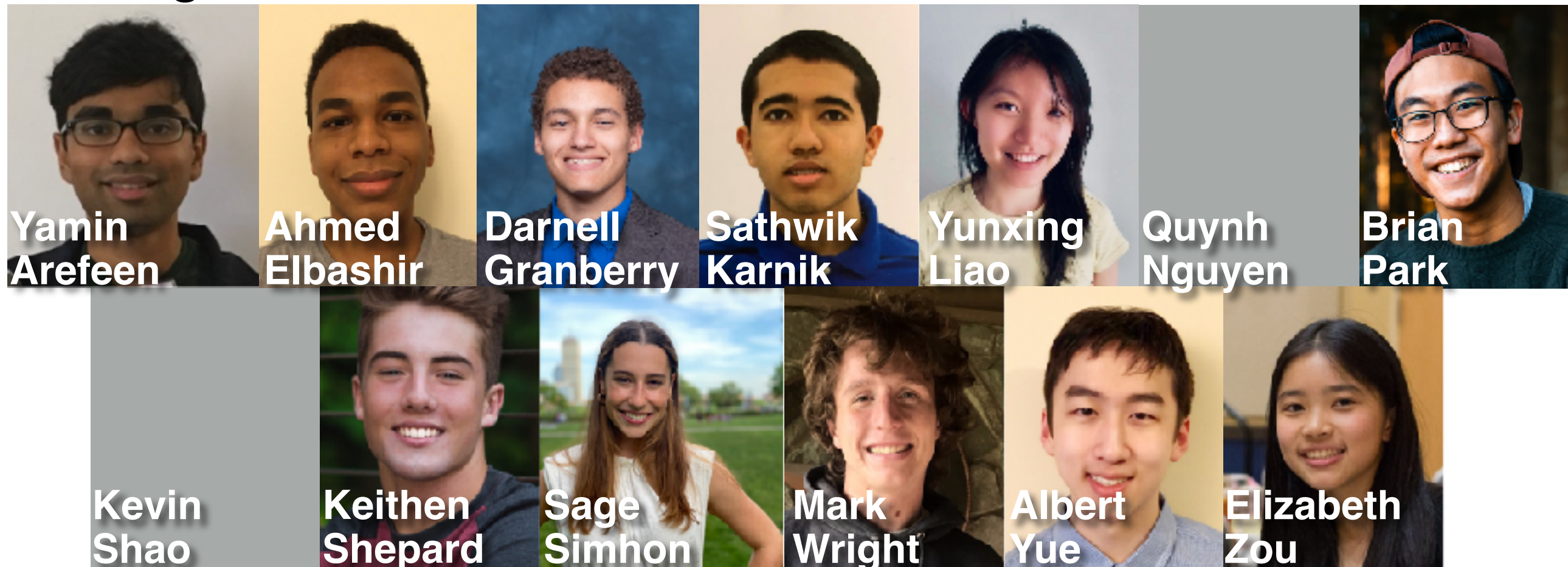


# 6.036: Introduction to Machine Learning, Staff

Instructors:



Teaching Assistants:



And Lab Assistants!

# Machine learning (ML): why & what

# Machine learning (ML): why & what

**nature**

Search Login

Content ▾ About ▾ Publish ▾

NEWS | 22 July 2021

## **DeepMind's AI predicts structures for a vast trove of proteins**

AlphaFold neural network produced a 'totally transformative' database of more than 350,000 structures from *Homo sapiens* and 20 model organisms.

Ewen Callaway

Underpinning the latest version of AlphaFold is a novel machine learning approach that incorporates physical and biological knowledge about protein structure, leveraging multi-sequence alignments, into the design of the deep learning algorithm.



# Machine learning (ML): why & what

nature

Search Login

TechRepublic.

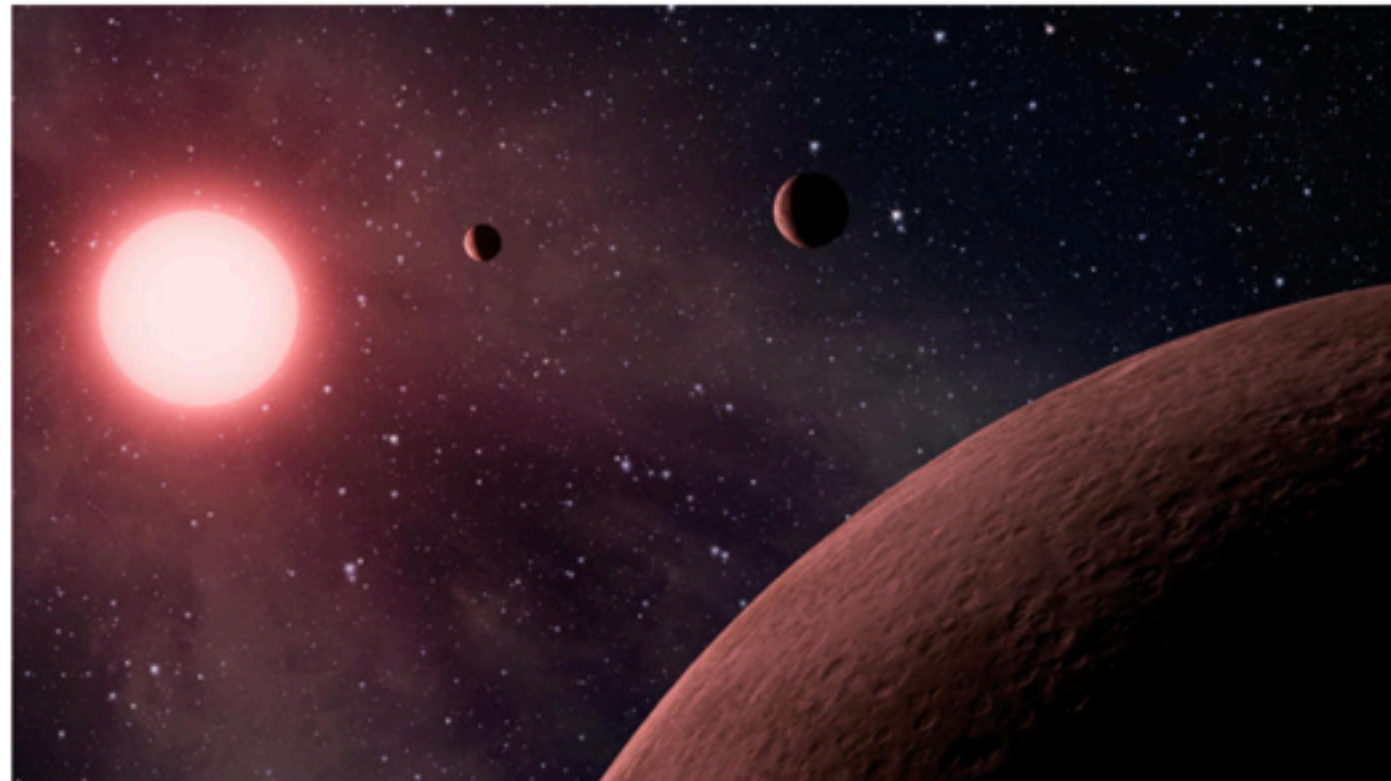


## Machine learning algorithm confirms 50 new exoplanets in historic first



by **R. Dallan Adams** in **Innovation** on August 26, 2020, 9:07 AM PST

A new machine learning technique can be used to sift through massive datasets to discern exoplanets from false positives.





# Machine learning (ML): why & what

nature

Search Login

TechRepublic.



TECHNOLOGY

JULY 22, 2021 / 9:23 AM / UPDATED 2 MONTHS AGO

## FEATURE-Catching fire: AI is helping scarce firefighters better predict blazes

By Avi Asher-Schapiro, Thomson Reuters Foundation

One of the most complex tools developed by researchers in recent years is the Potential Control Locations (PCL) algorithm, which uses machine learning to suggest where firefighters should place their control lines during a blaze.

# Machine learning (ML): why & what

nature

Search Login

TechRepublic

REUTERS | The Echo Chamber

TECHNOLOGY

JULY 22, 2021

*A small group of lawyers and its outsized influence at the U.S. Supreme Court*

Machine

FEATU

scarce

By Avi Ashe



TOP TIER: In handling appeals heard by the U.S. Supreme Court, 75 lawyers have stood out – most for their success at getting cases before the high court, others for how often they argue those cases, and some for both reasons. Most of the 75 work at law firms that primarily represent businesses.

## At America's court of last resort, a handful of lawyers now dominates the docket

By [Joan Biskupic](#), [Janet Roberts](#) and [John Shiffman](#)

Conten

NEWS

De

str

of

Alpha

datab

sapie

Ewen



# Machine learning (ML): why & what

**nature**

Conten

TechRepublic

REUTE

TECHNOLOGY

JULY 22, 2021

A small g

NEWS

Mach new e

De str of

Alpha datab sapie

Ewen

FEATU scarce

One dev rec

Cor

algo lear fire con

At A resc now

TOP TIER: In success at ge both reasons

By Av Ash

A new mach massive dat

By Joan Bi

**npr** WAMU 88.5 AMERICAN UNIVERSITY RADIO

DONATE

TECHNOLOGY

## Massachusetts Pioneers Rules For Police Use Of Facial Recognition Tech

May 7, 2021 · 6:00 AM ET

EMMA PEASLEE

Surveillance cameras, like the one here in Boston, are used throughout Massachusetts. The state now regulates how police use facial recognition technology.

[<https://www.npr.org/2021/05/07/982709480/massachusetts-pioneers-rules-for-police-use-of-facial-recognition-tech>]

# Machine learning (ML): why & what

**nature**

Conten

TechRepublic

REUTERS

TECHNOLOGY

Mach new e

FEATU scarce

One dev

reco

Cor

algo

lear

fire

con

At A resc now

By Joan Bi

**npr** WAMU 88.5

DONATE

TECHNOLOGY

Aug 13, 2021, 01:51pm EDT | 2,064 views

Mas For Rec

May 7, 20

## Upstart: Can AI Kill The FICO Score?

Tom Taulli  
Contributor  
Entrepreneurs  
I write about tech & finance.

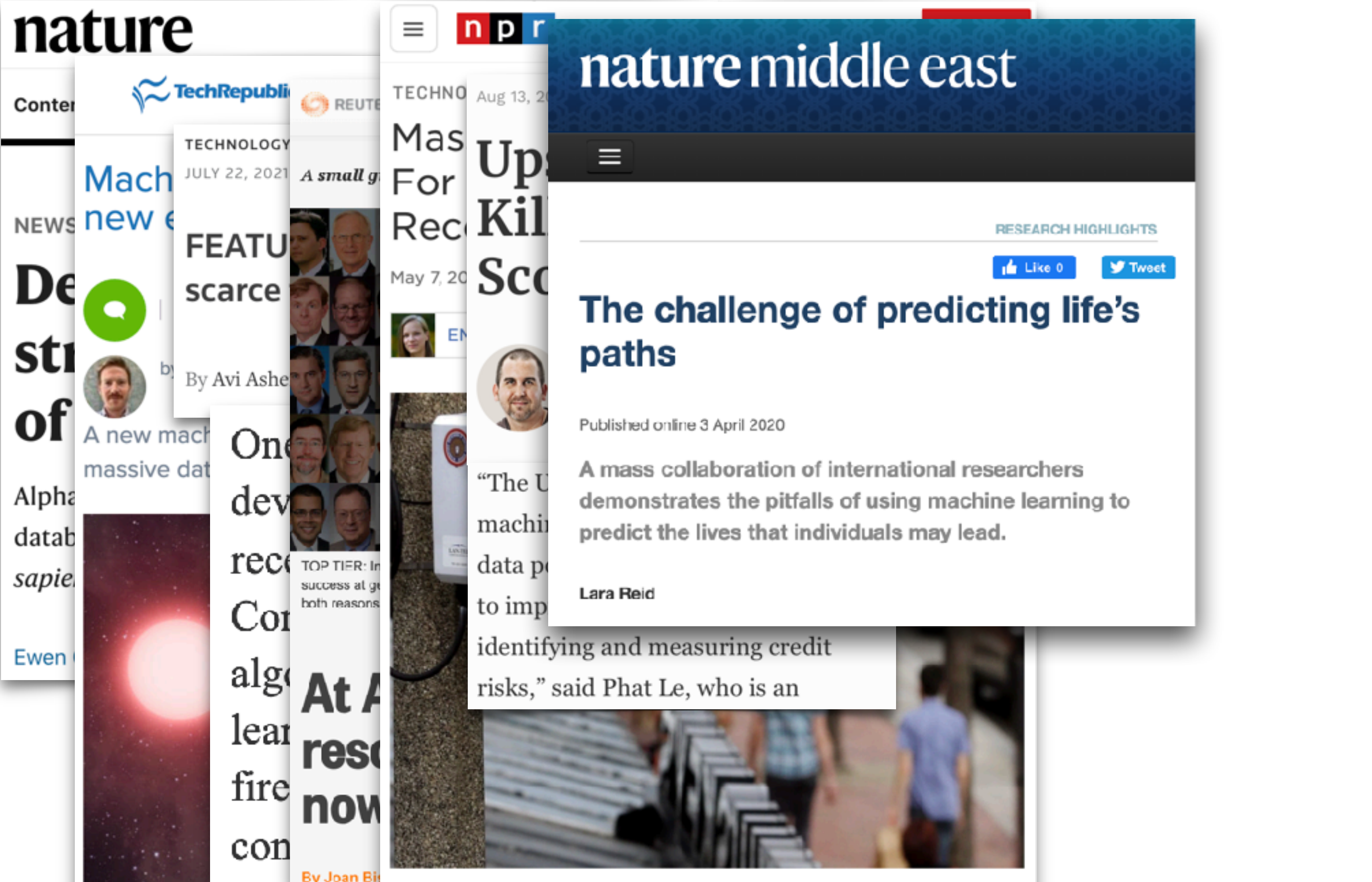
Follow

“The Upstart system uses AI and machine learning models with 1,600 data points and 15 billion cells of data to improve accuracy in terms of identifying and measuring credit risks,” said Phat Le, who is an

Surveillance cameras, like the one here in Boston. are used throughout Massachusetts. The state now regulates how p

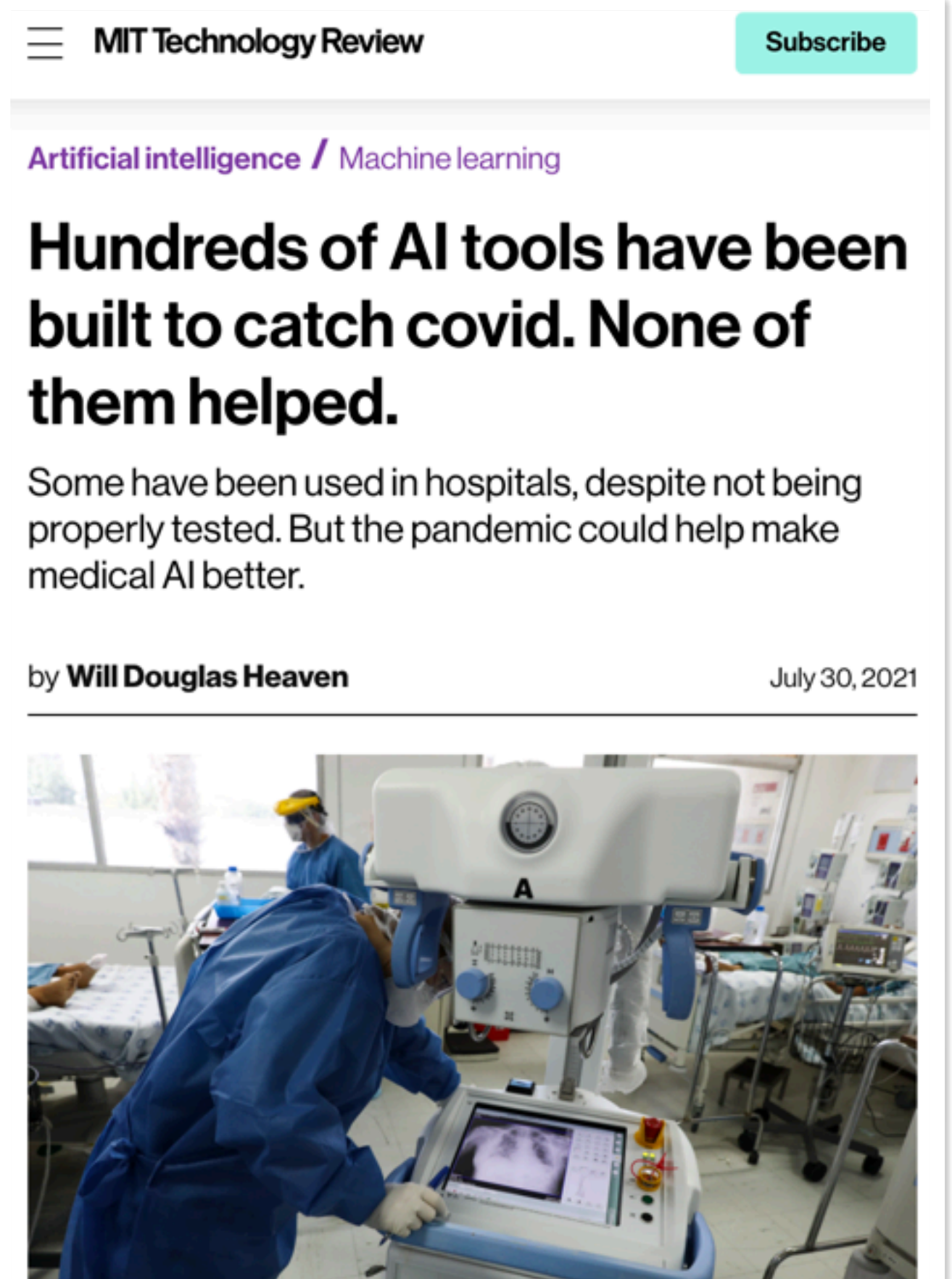
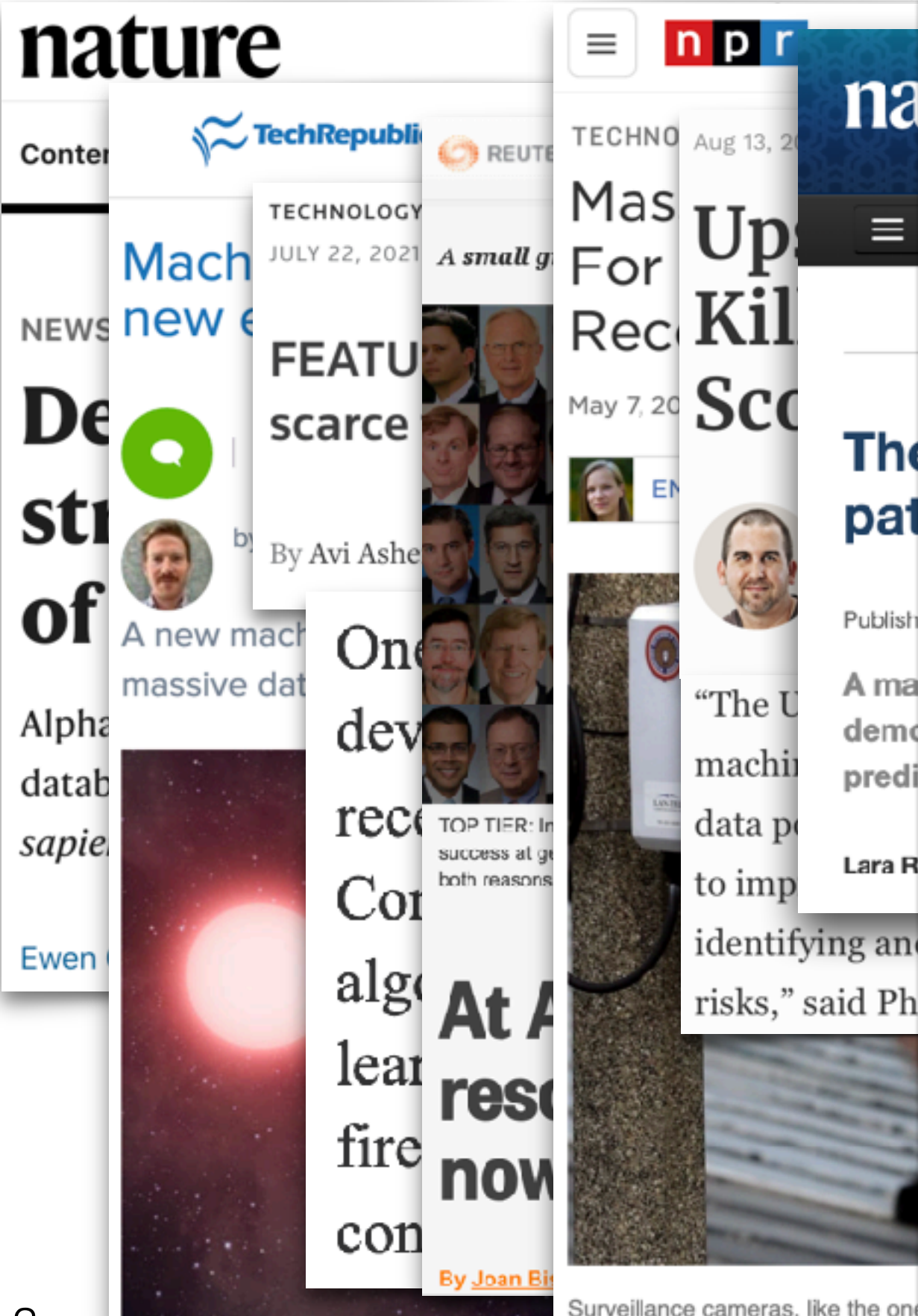


# Machine learning (ML): why & what





# Machine learning (ML): why & what

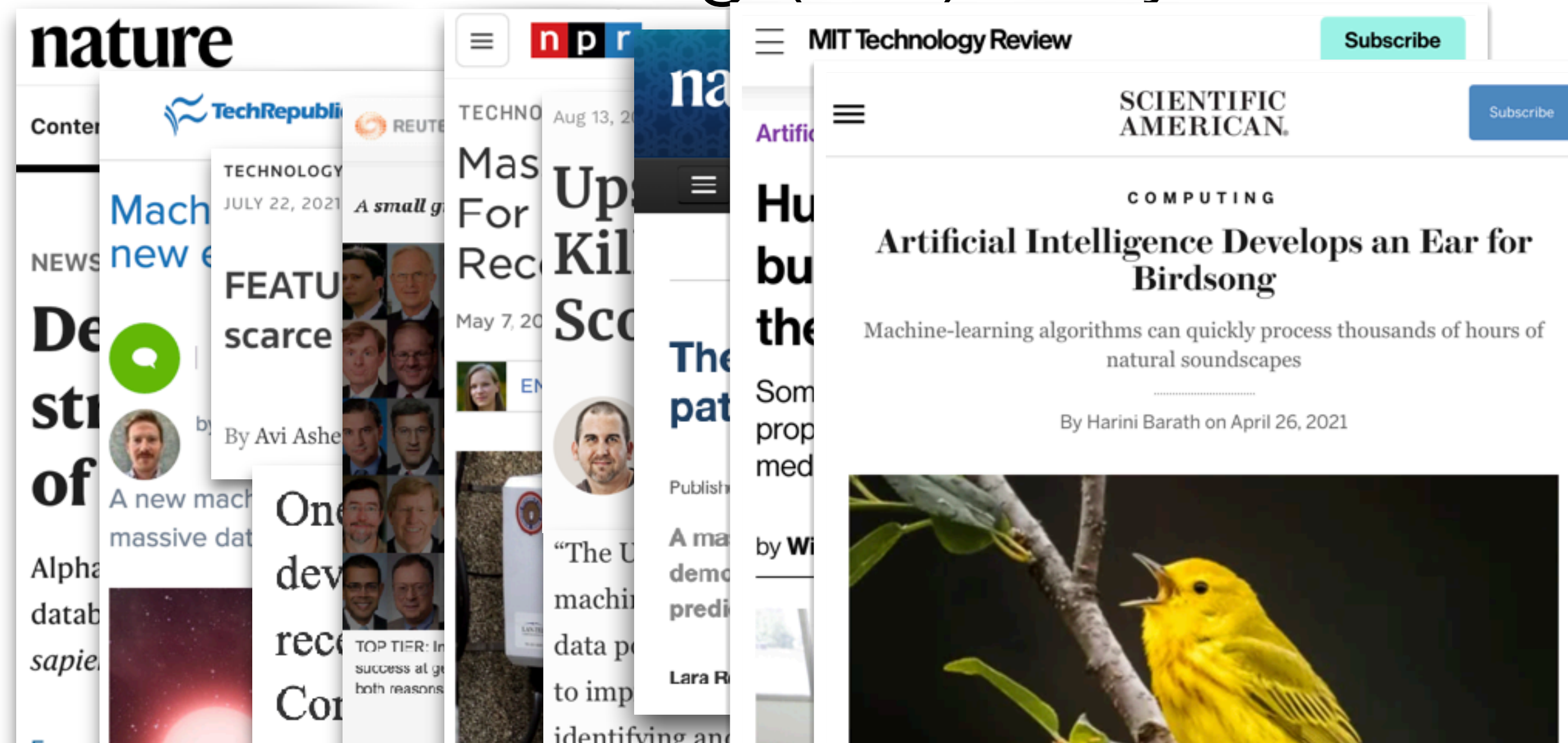




# Machine learning (ML): why & what

The image is a collage of various news articles and a featured article. On the left, there are several overlapping article snippets from sources like Nature, TechRepublic, and NPR, with headlines such as 'Mach new e', 'De str of', 'Alpha datab sapie', 'Ewen', 'FEATU scarce', 'One dev', 'reco', 'Cor', 'algo', 'lear', 'fire', 'con', 'Mas For Rec', 'Up Kil Sco', 'The pat', 'Artific', 'Hu bu the', 'Som prop med', 'A ma demo predi', 'Lara R', 'TOP TIER: In success at ge both reasons', 'At A resc now', 'By Joan Bi', 'Surveillance cameras. like the one', 'The state now'. On the right, a full article from Scientific American is displayed. The article is titled 'Artificial Intelligence Develops an Ear for Birdsong' and is categorized under 'COMPUTING'. The sub-headline reads 'Machine-learning algorithms can quickly process thousands of hours of natural soundscapes'. The author is Harini Barath, and the date is April 26, 2021. Below the text is a photograph of a male yellow warbler perched on a branch, singing. The caption for the photo is 'Male yellow warbler in Yosemite National Park. Credit: Alice Cahill Getty Images'. At the top right of the Scientific American article, there is a 'Subscribe' button. At the top left of the MIT Technology Review article, there is a 'Subscribe' button.

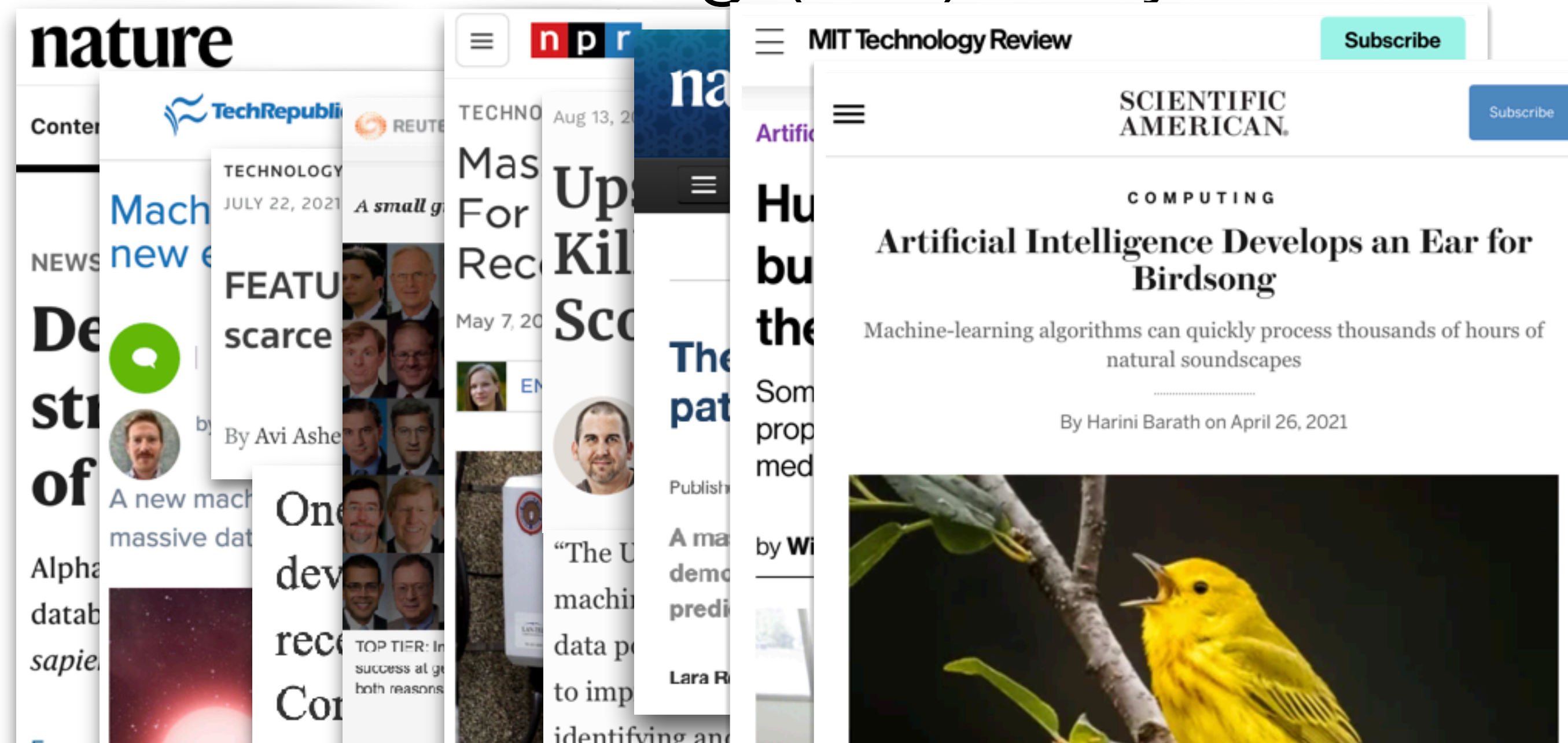
# Machine learning (ML): why & what



- **What is ML?**

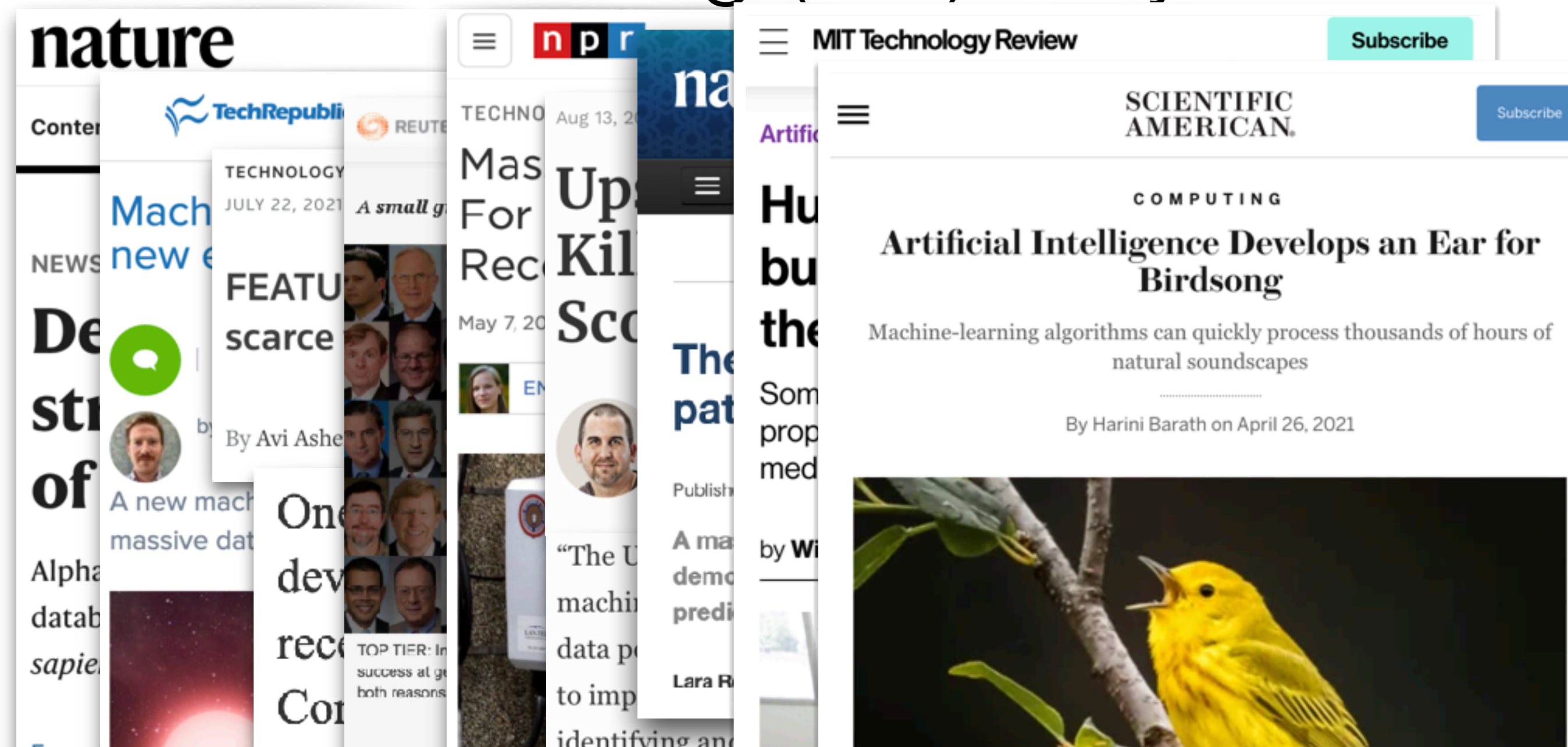


# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)

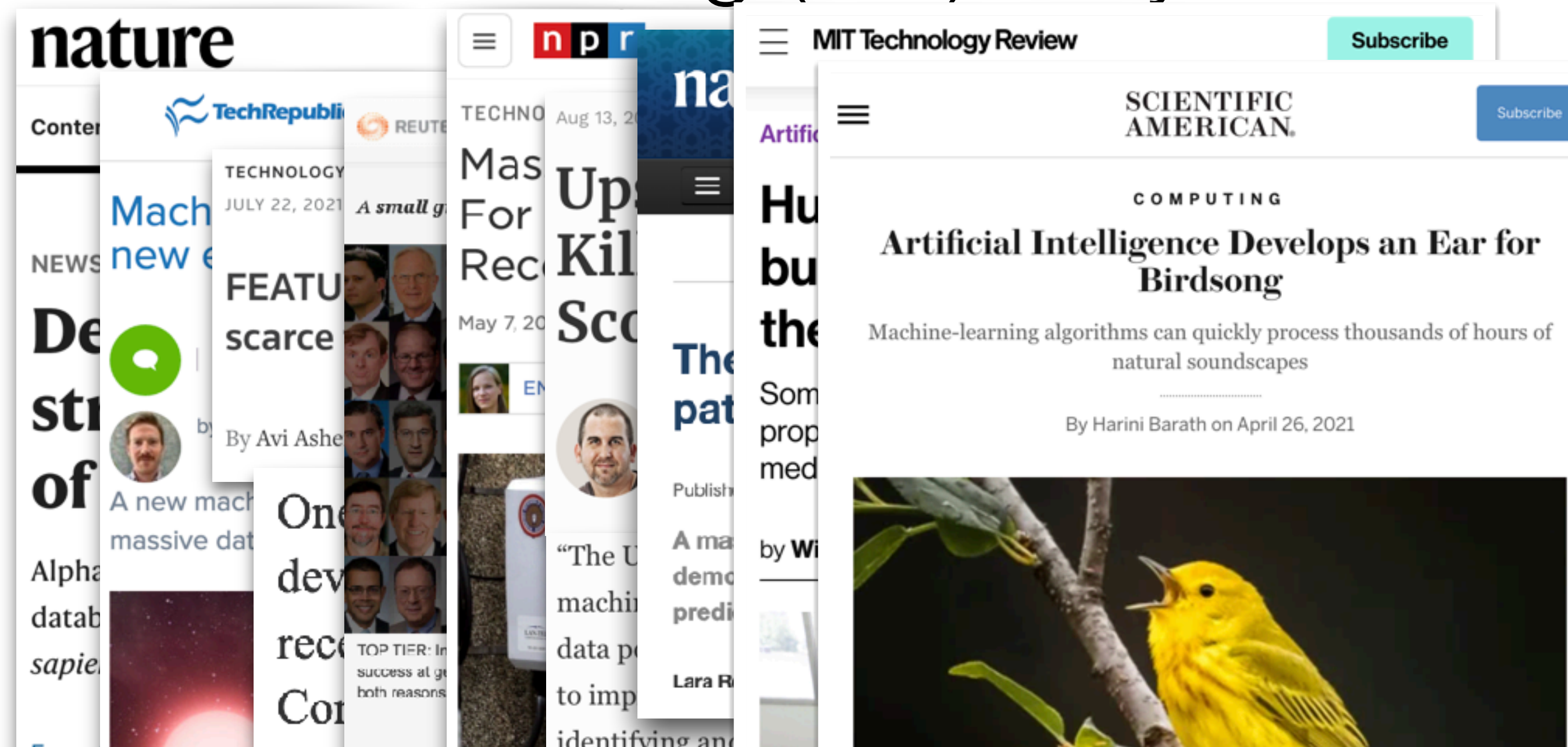
# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate



# Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate
- **Notes:** ML is a tool with pros & cons. ML is built on math

# Getting started: regression



# Getting started: regression

**Example:** predict pollution level

# Getting started: regression

**Example:** predict pollution level

**What do we have?**

# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points

# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$

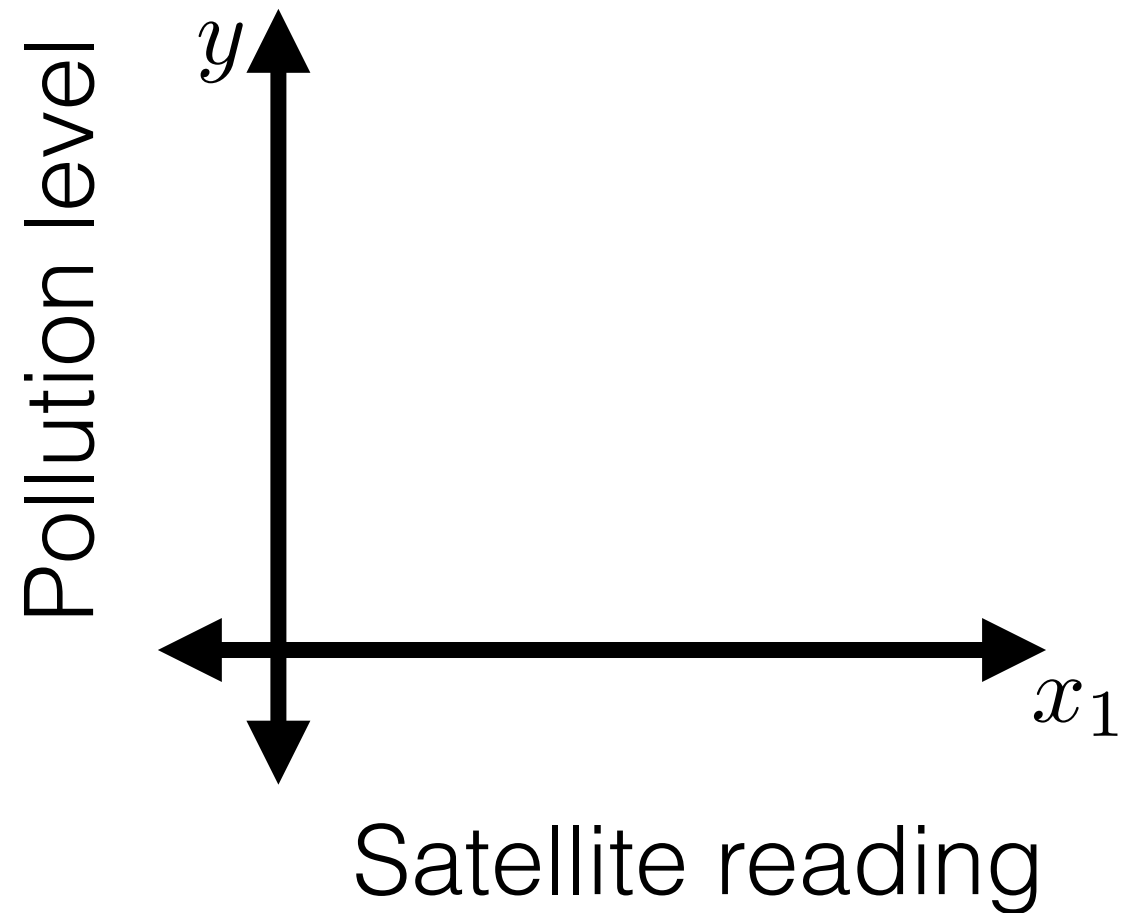


# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$

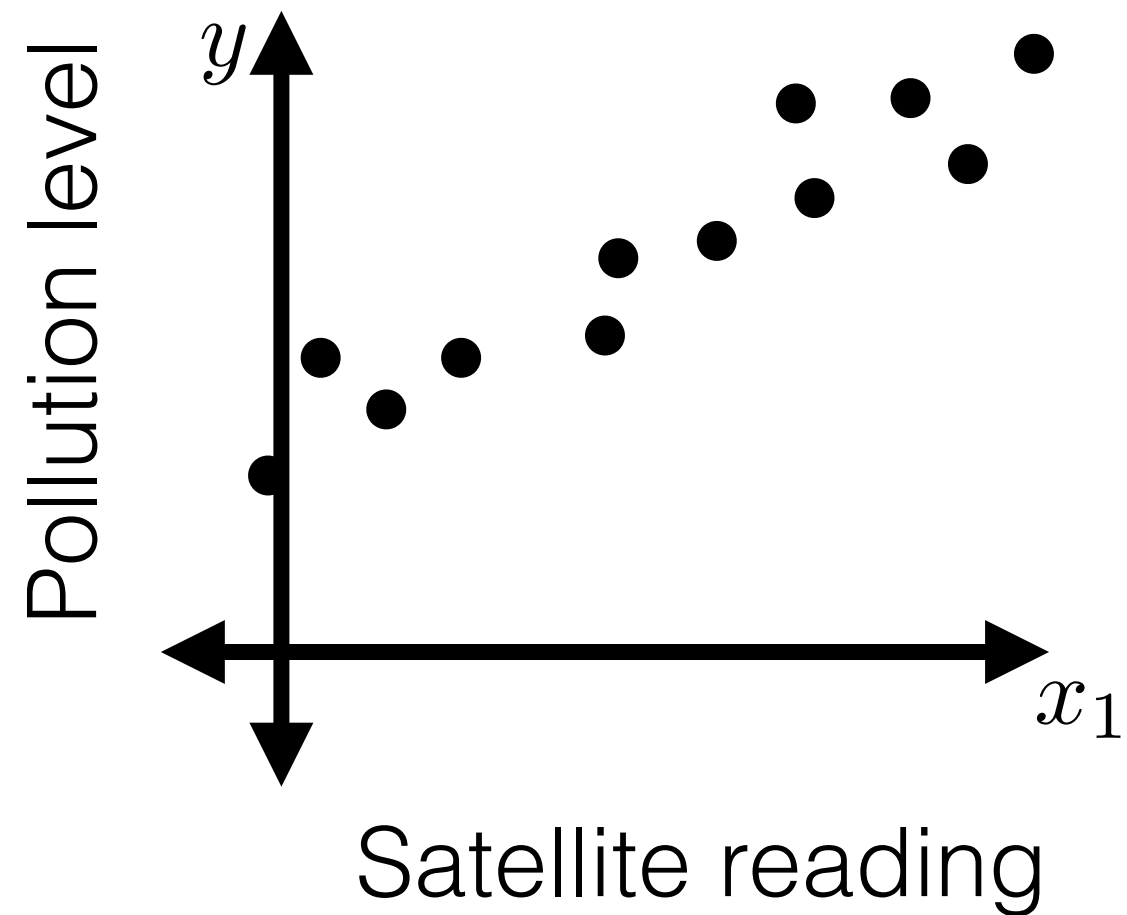


# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$

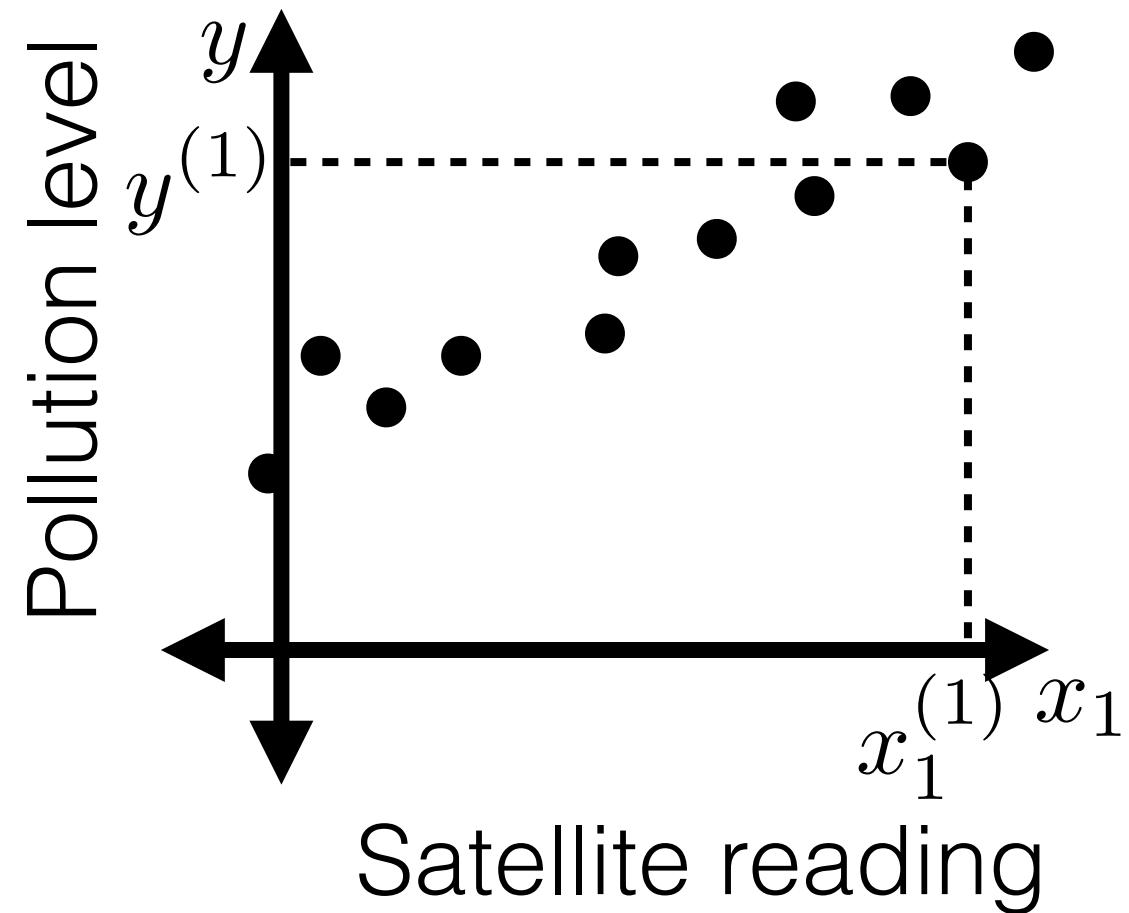


# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$

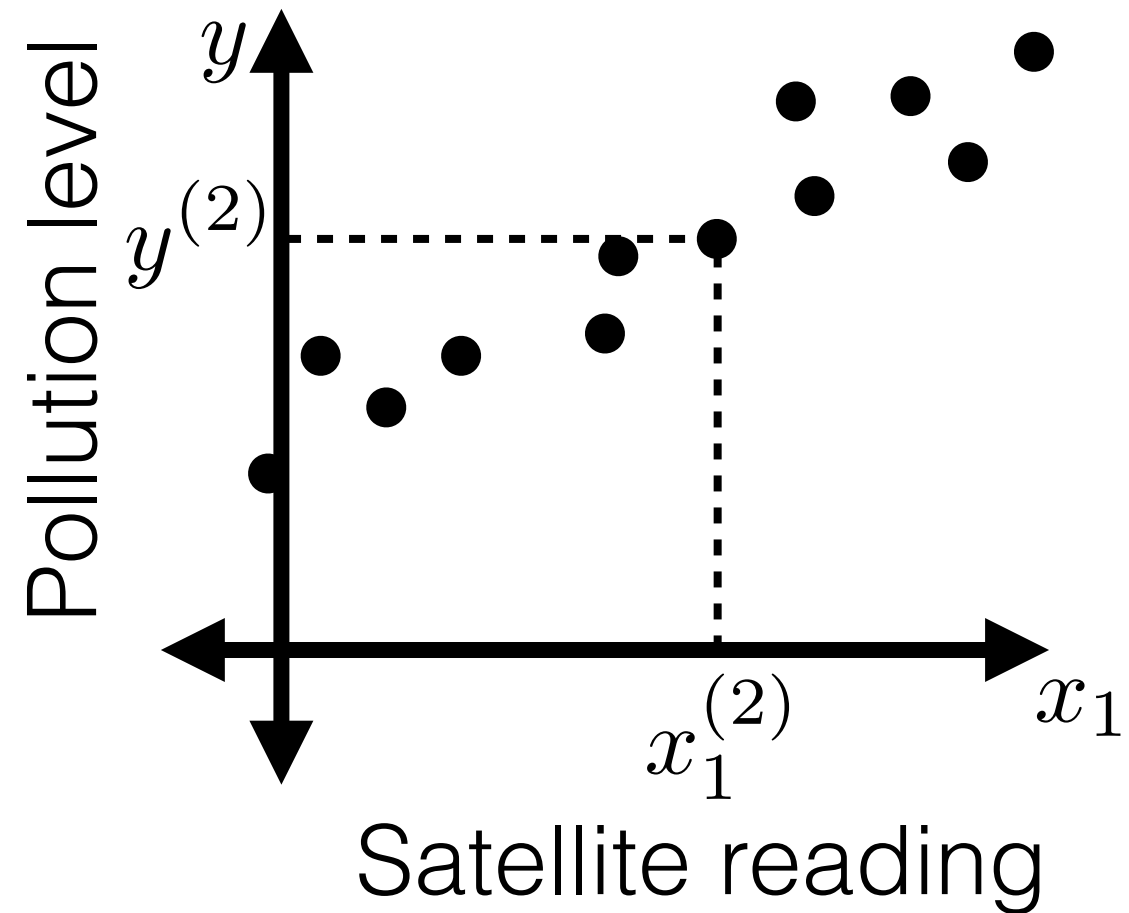


# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

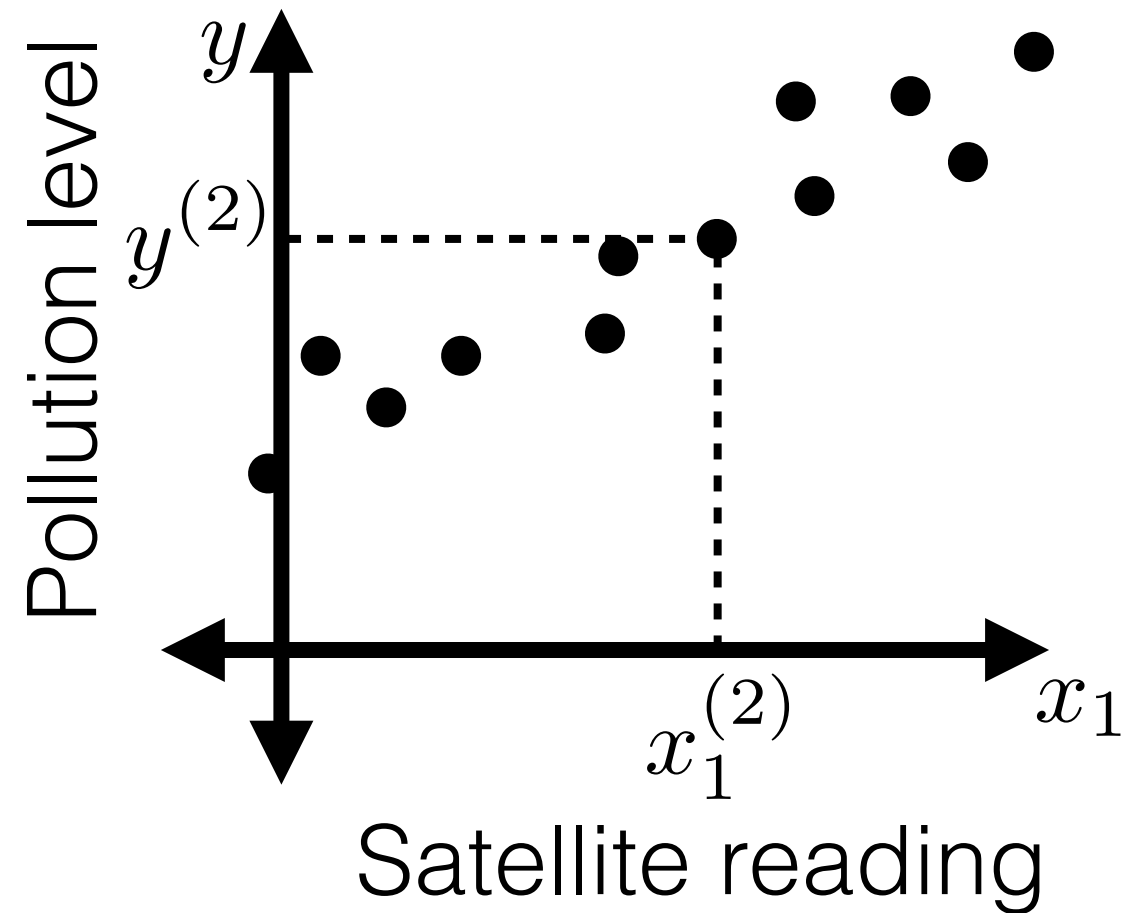
- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

- Feature vector

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$

- Label  $y^{(i)} \in \mathbb{R}$

- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

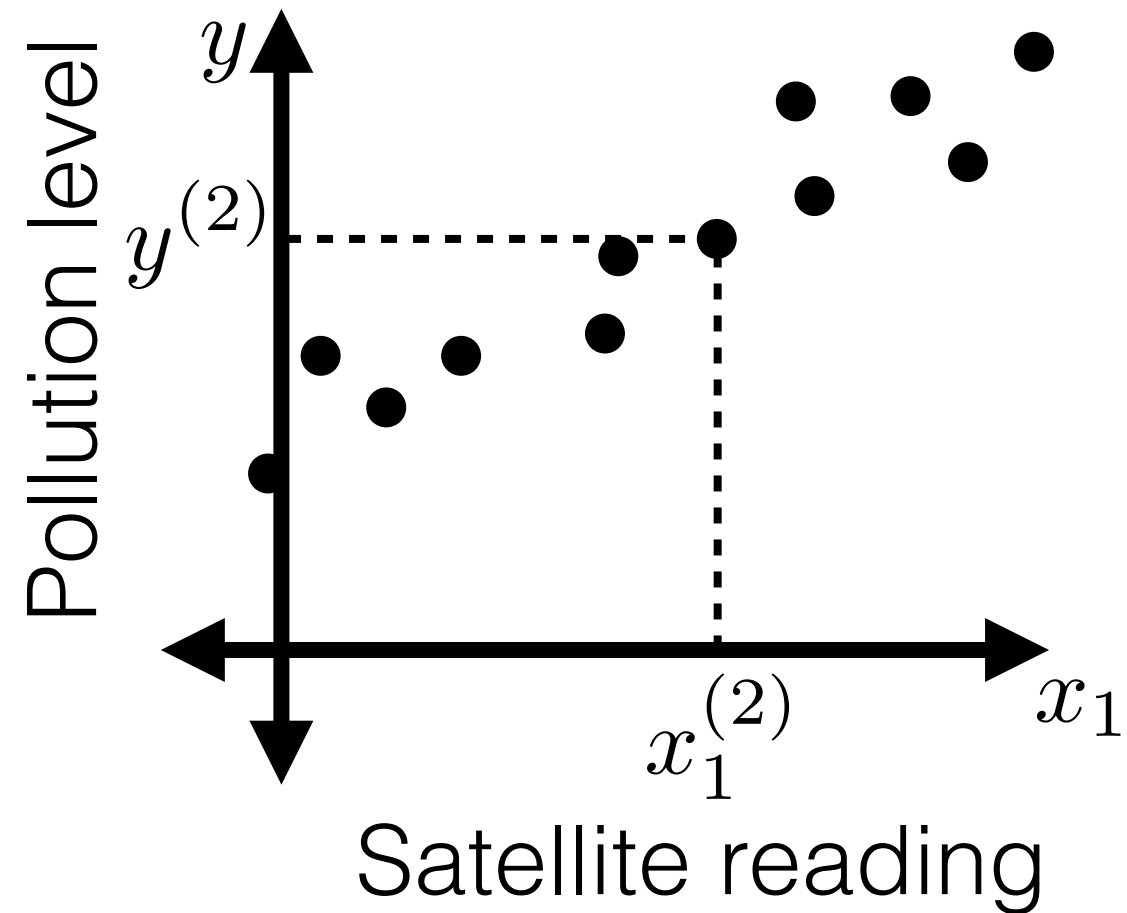
- Feature vector

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$

- Label  $y^{(i)} \in \mathbb{R}$

- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

**What do we want?**





# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

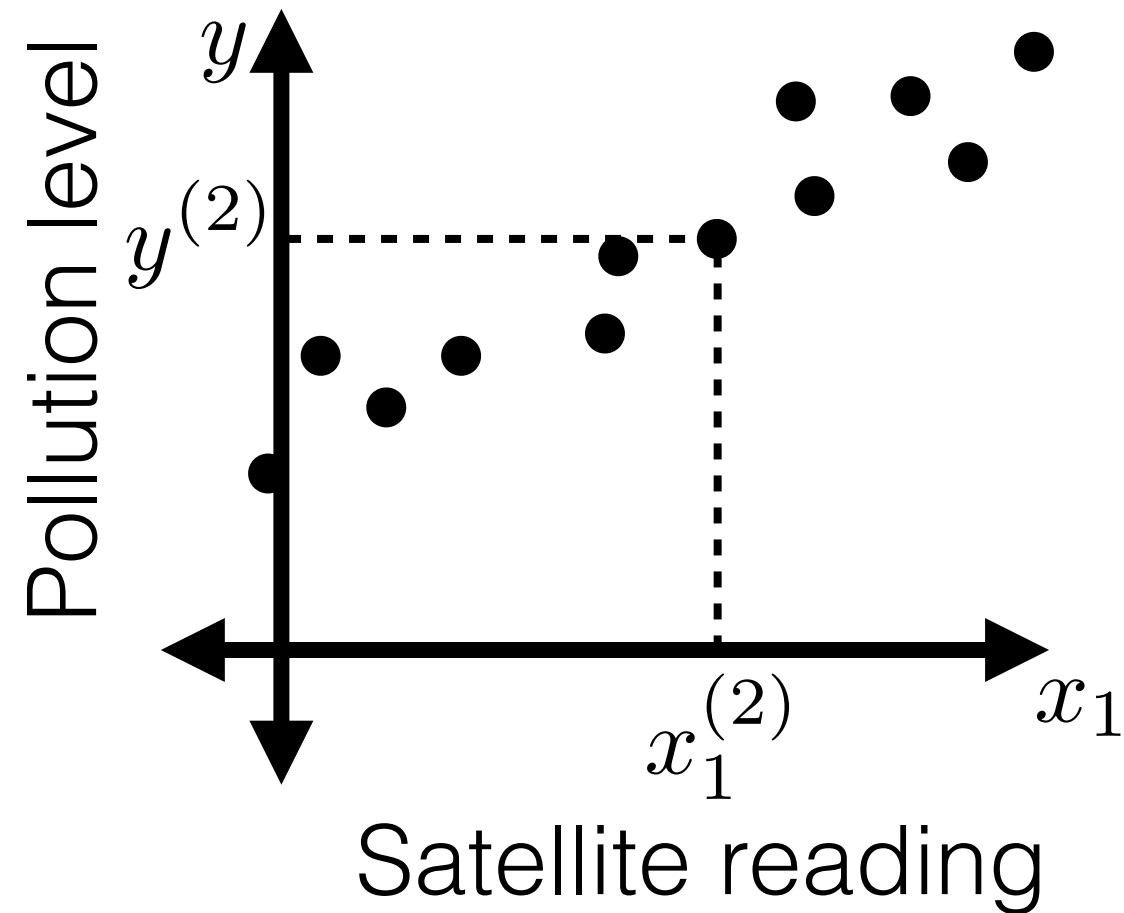
- Feature vector

$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$

- Label  $y^{(i)} \in \mathbb{R}$

- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

**What do we want?** A good way to label new points



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

- Feature vector

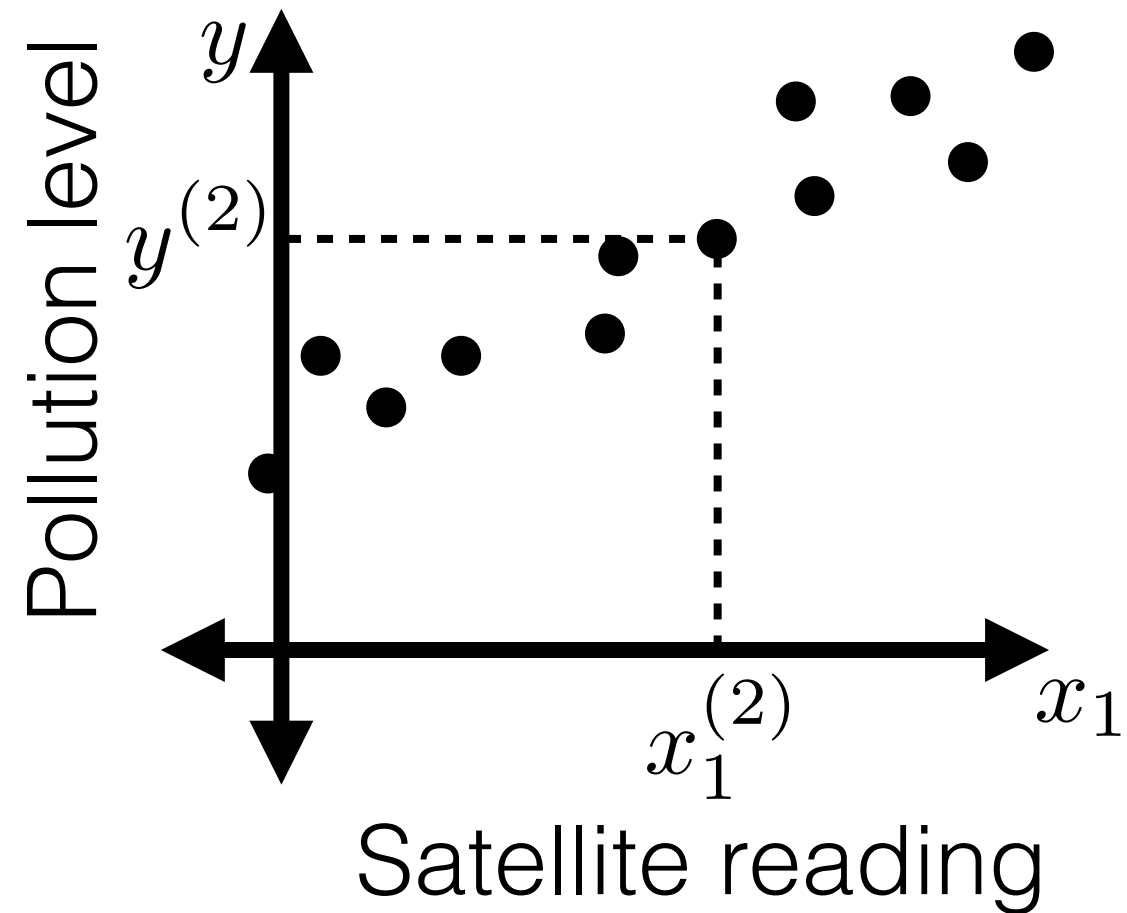
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$

- Label  $y^{(i)} \in \mathbb{R}$

- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

**What do we want?** A good way to label new points

- How to label?



# Getting started: regression

**Example:** predict pollution level

**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$

- Feature vector

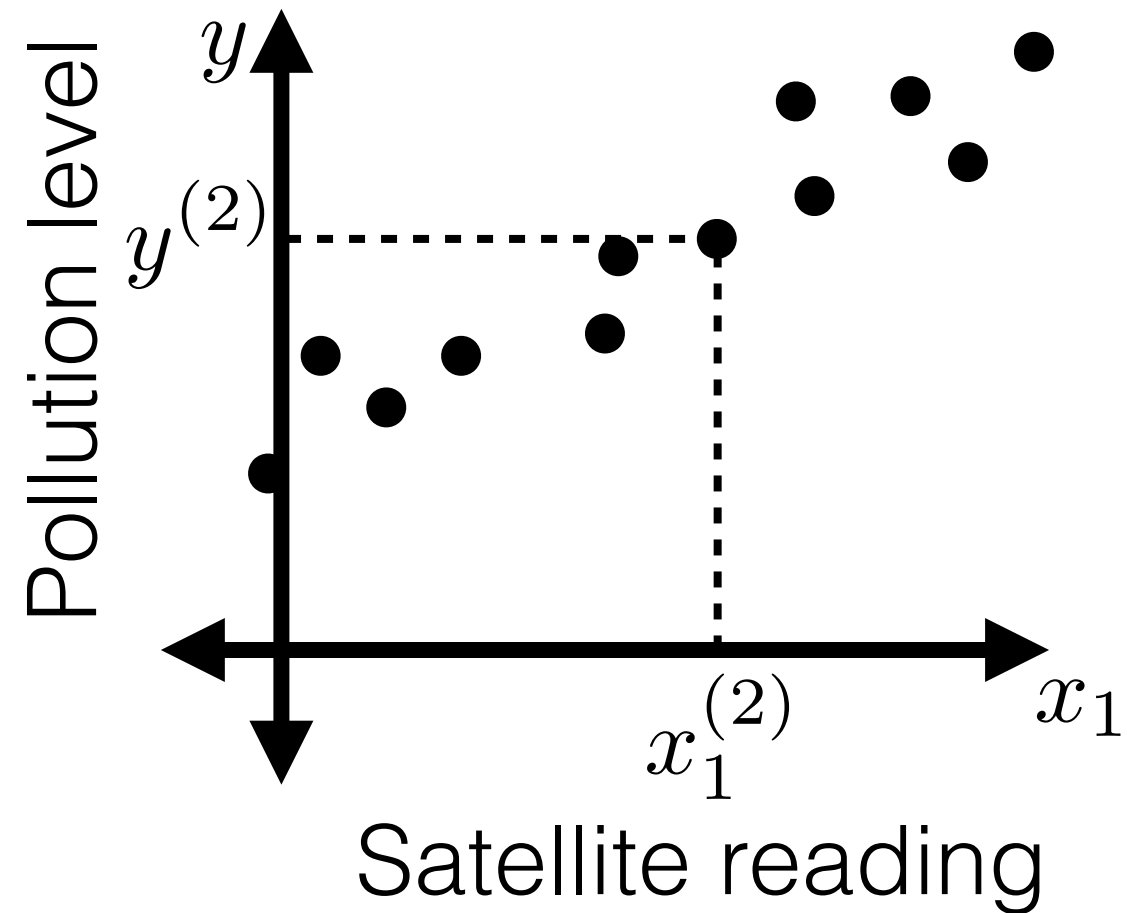
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$

- Label  $y^{(i)} \in \mathbb{R}$

- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$

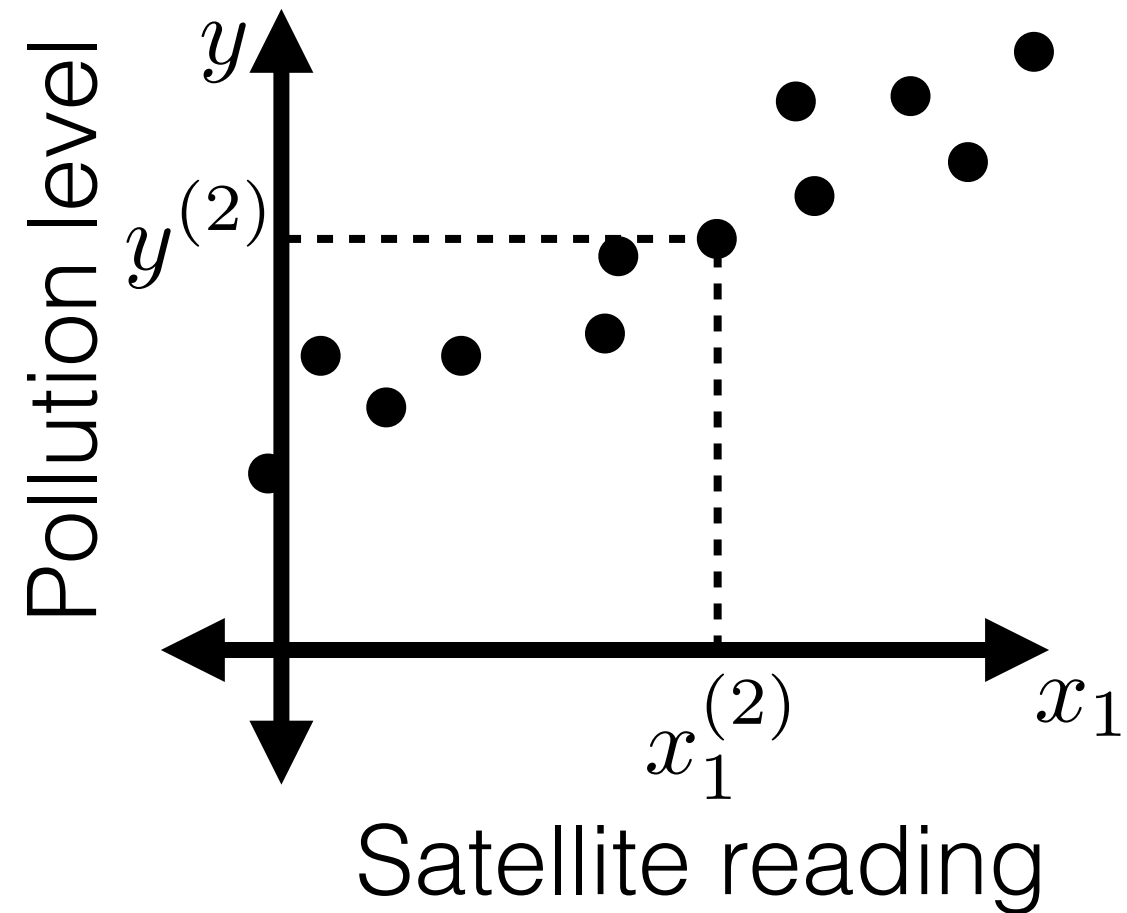


# Getting started: regression

**Example:** predict pollution level

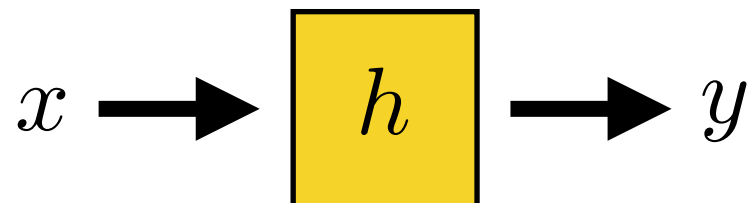
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$

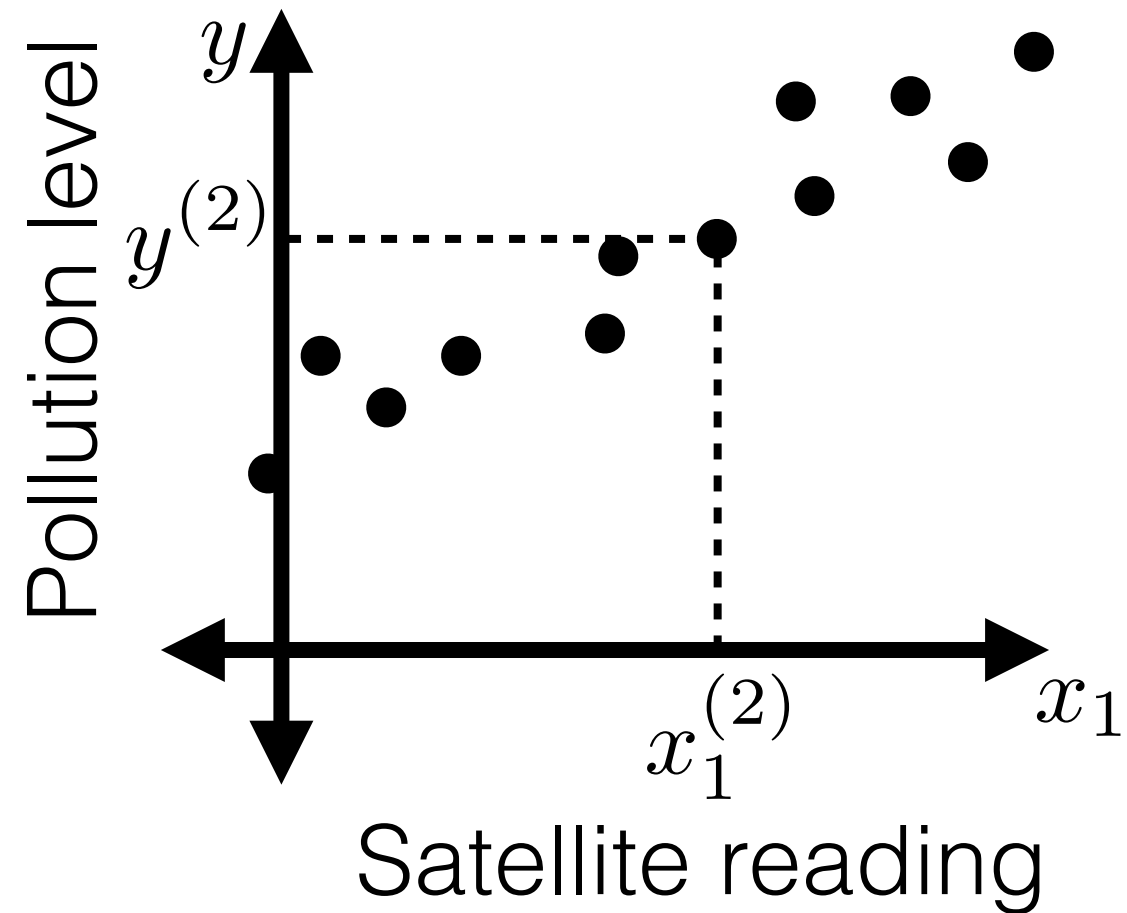


# Getting started: regression

**Example:** predict pollution level

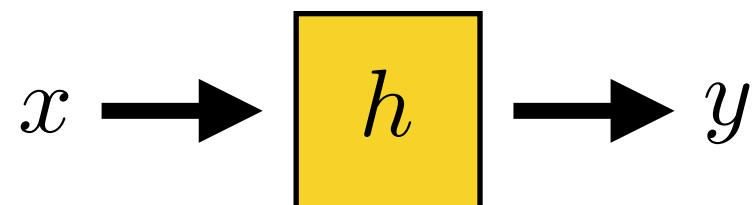
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$



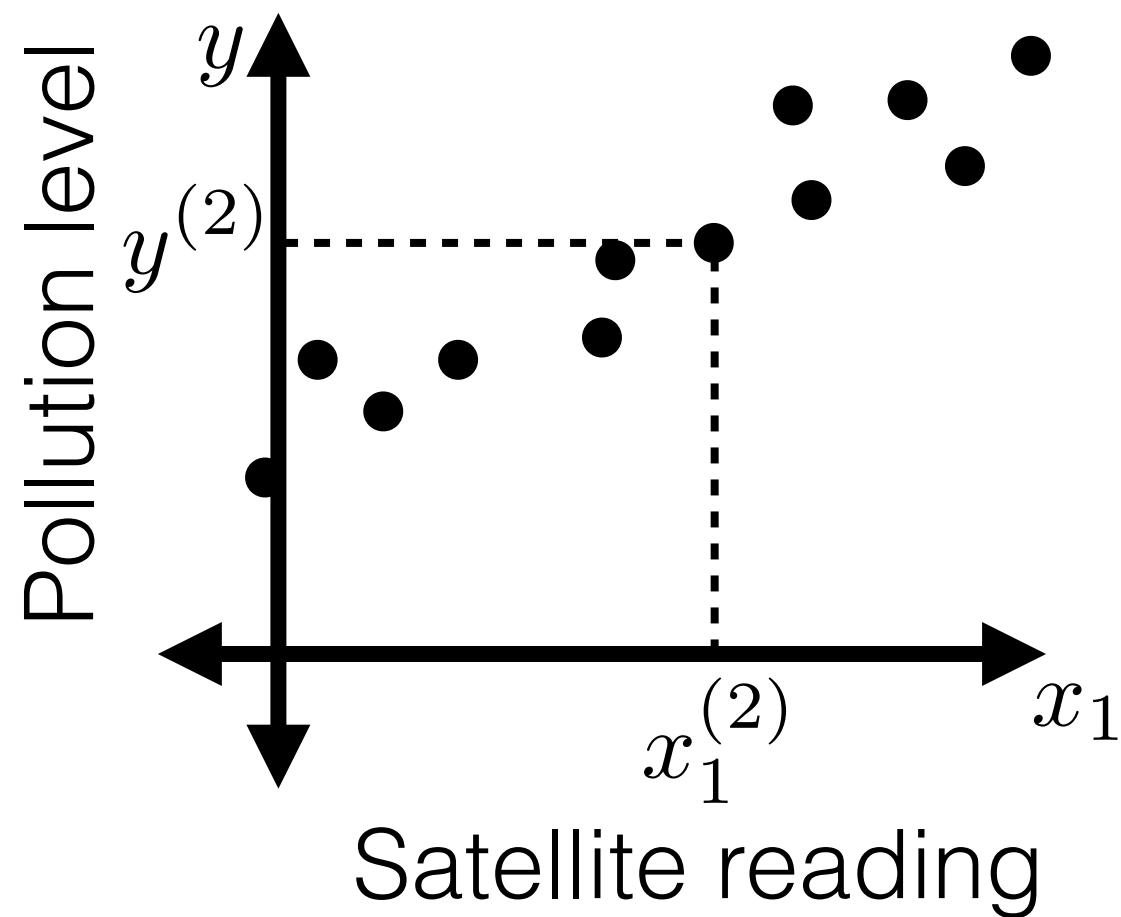
- Example  $h$ : For any  $x$ ,  $h(x) = 1,000,000$

# Getting started: regression

**Example:** predict pollution level

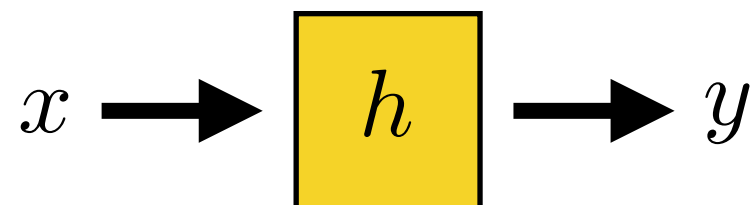
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$



Is this a hypothesis?

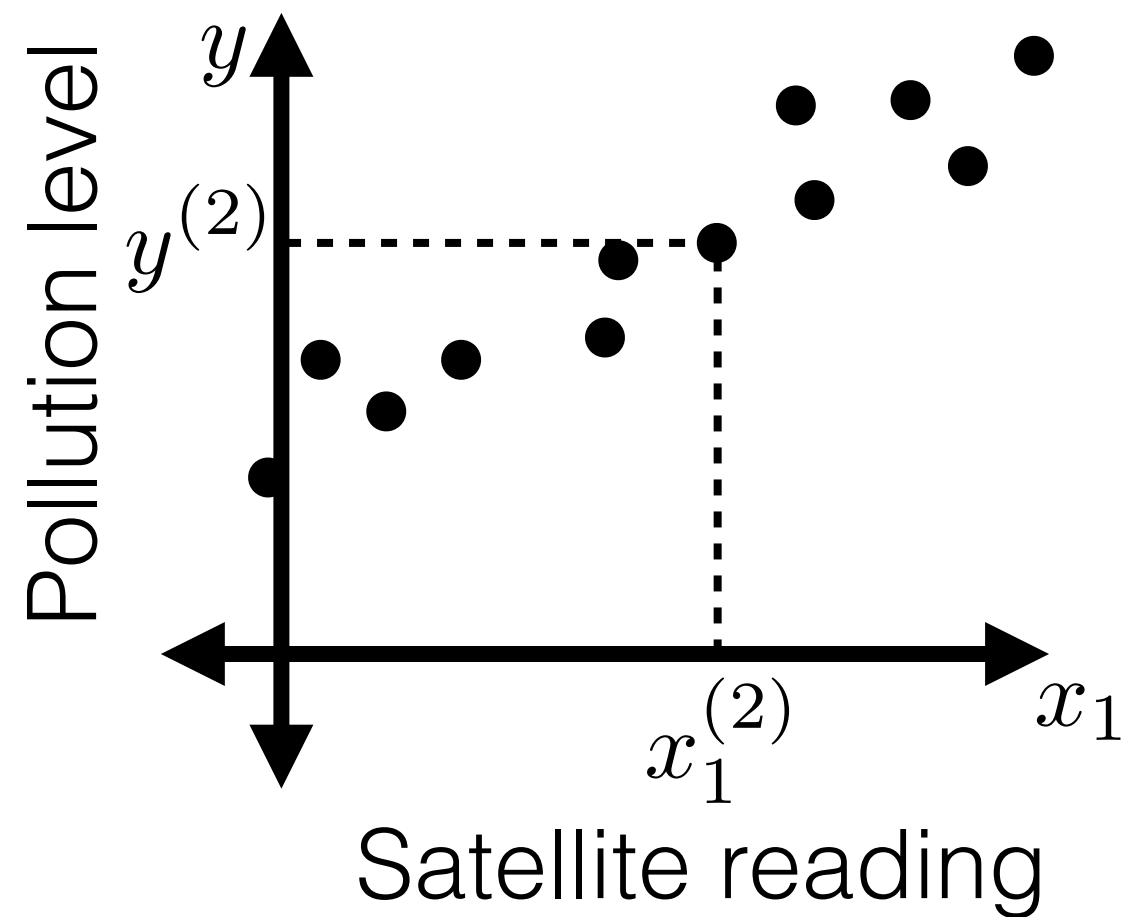
- Example  $h$ : For any  $x$ ,  $h(x) = 1,000,000$

# Getting started: regression

**Example:** predict pollution level

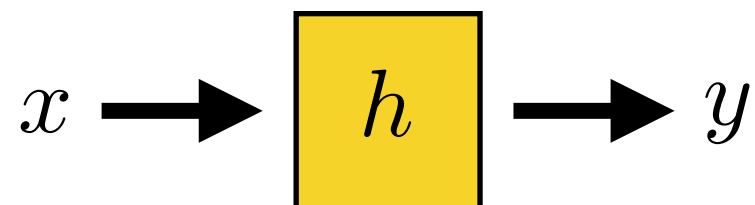
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$



Is this a **good** hypothesis?

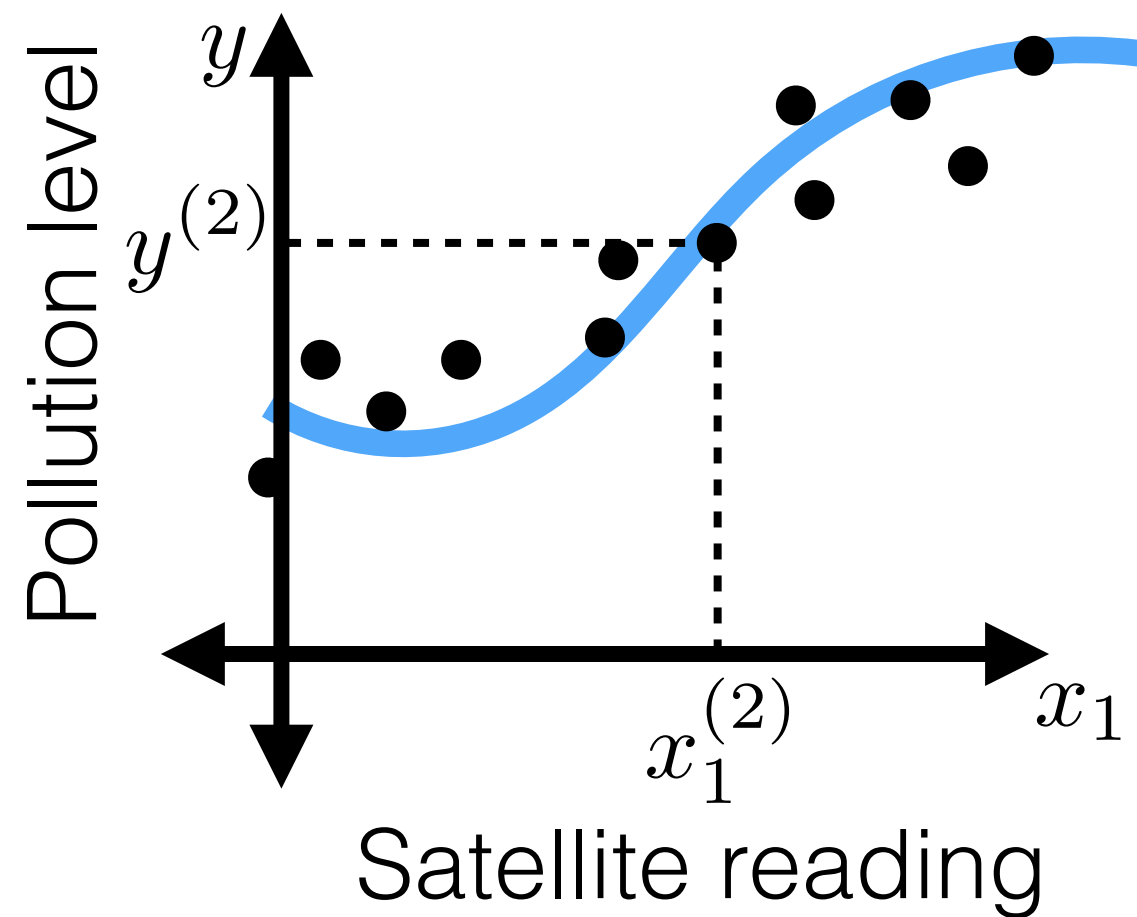
- Example  $h$ : For any  $x$ ,  $h(x) = 1,000,000$

# Getting started: regression

**Example:** predict pollution level

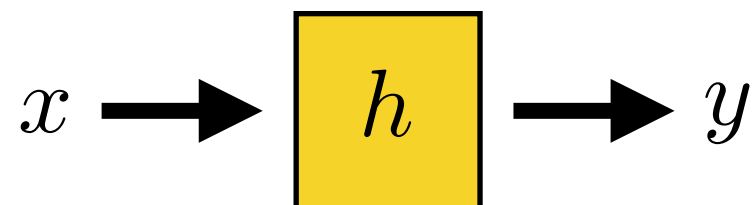
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$



Is this a **good** hypothesis?

- Example  $h$ : For any  $x$ ,  $h(x) = 1,000,000$

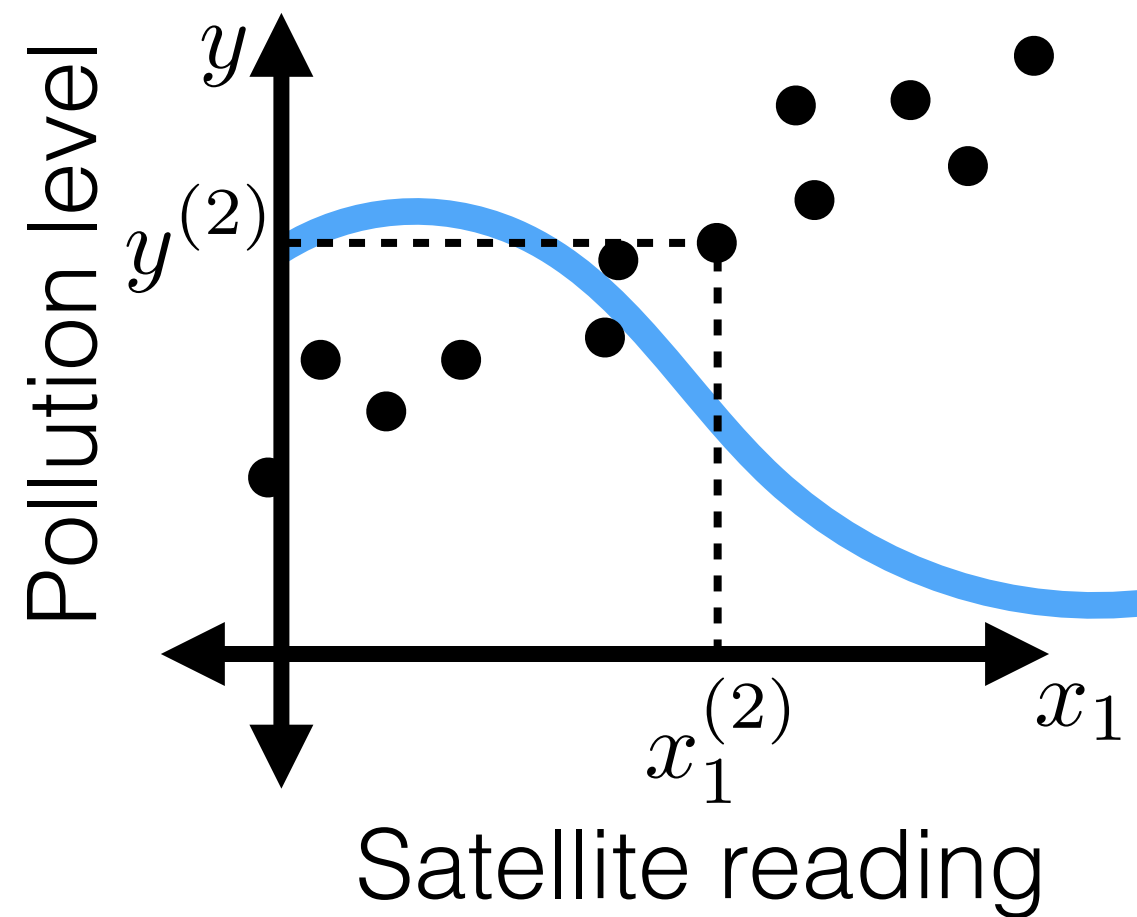


# Getting started: regression

**Example:** predict pollution level

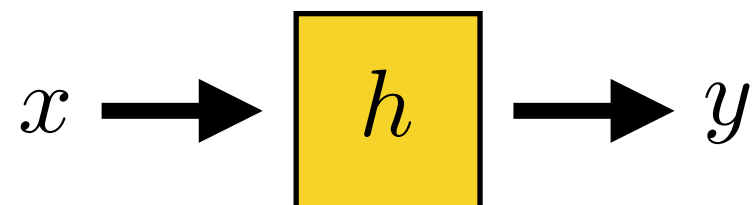
**What do we have?** (Training) data

- $n$  training data points
- For data point  $i \in \{1, \dots, n\}$ 
  - Feature vector  
$$x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$$
  - Label  $y^{(i)} \in \mathbb{R}$
- Training data  $\mathcal{D}_n = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$



**What do we want?** A good way to label new points

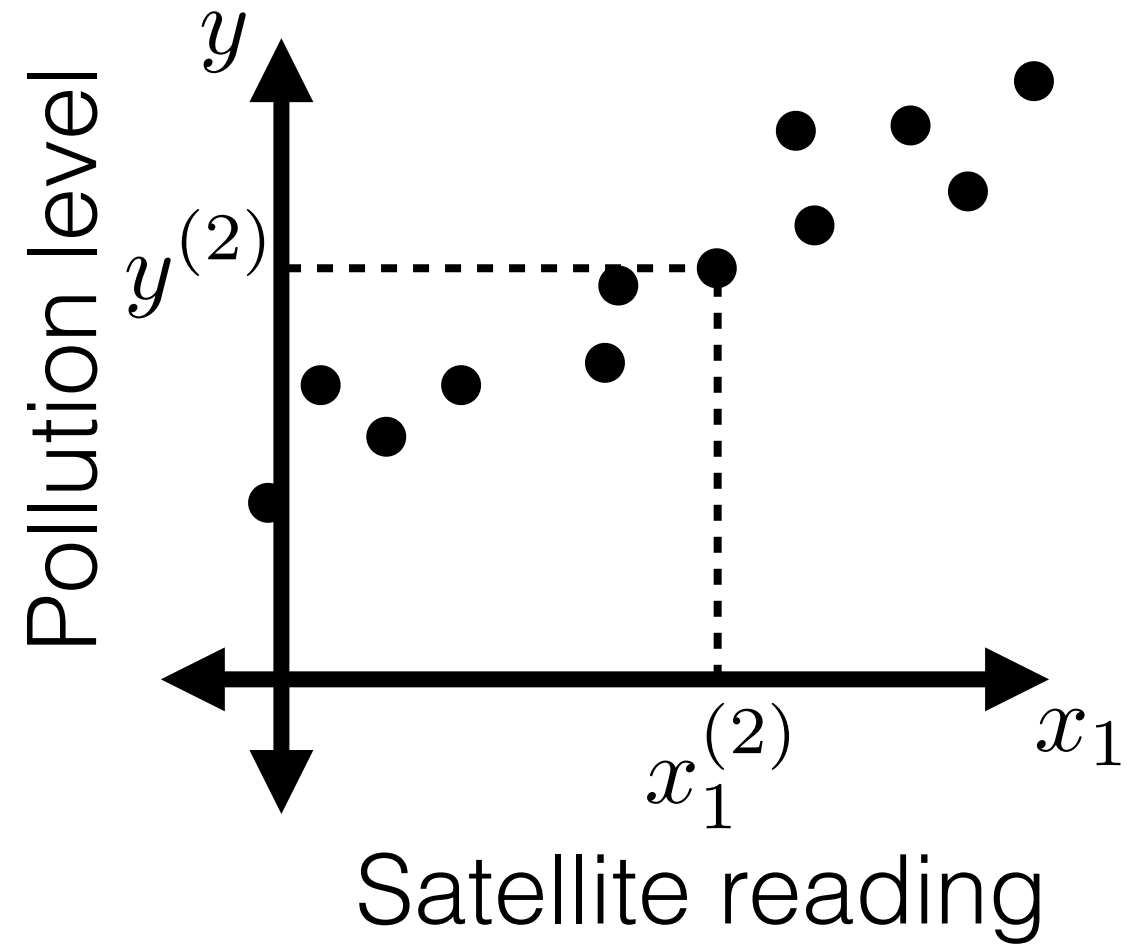
- How to label? Hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$



Is this a **good** hypothesis?

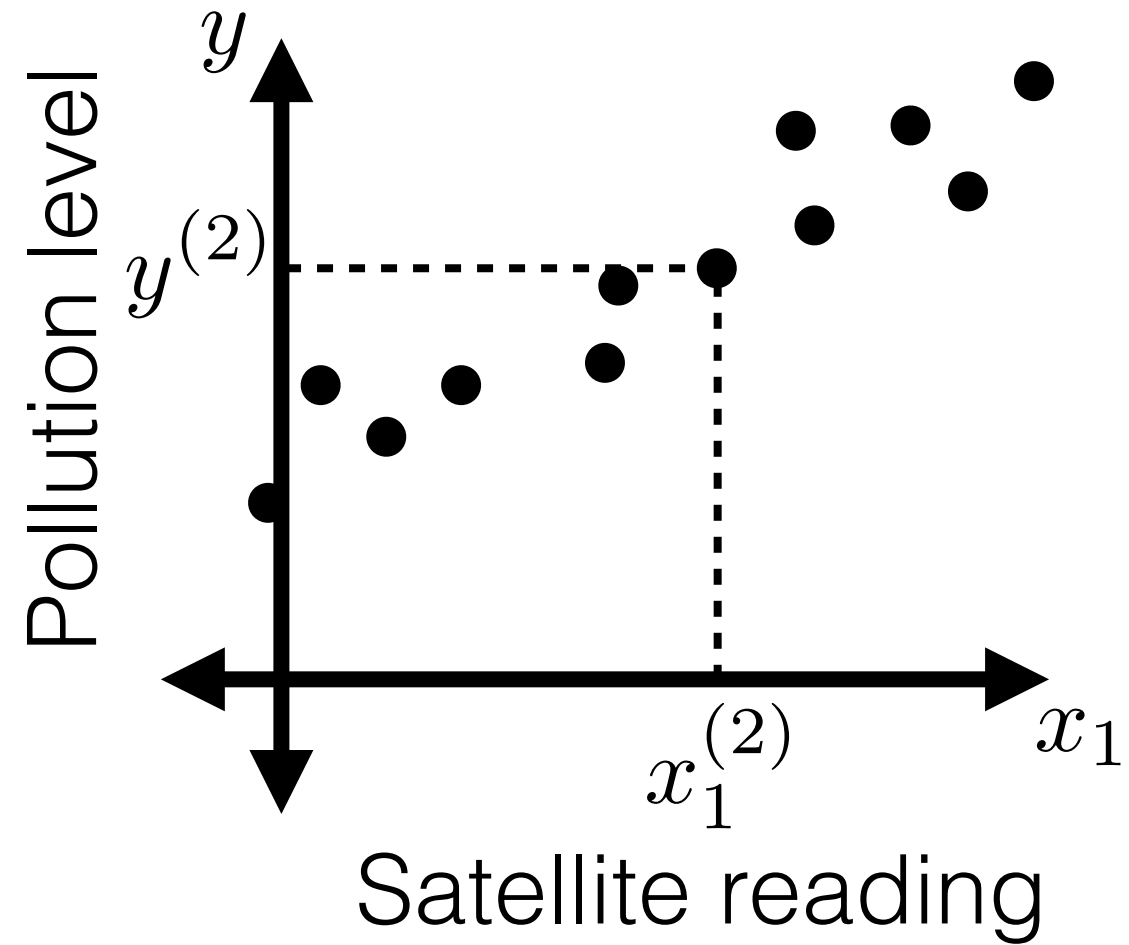
- Example  $h$ : For any  $x$ ,  $h(x) = 1,000,000$

# Linear regressors



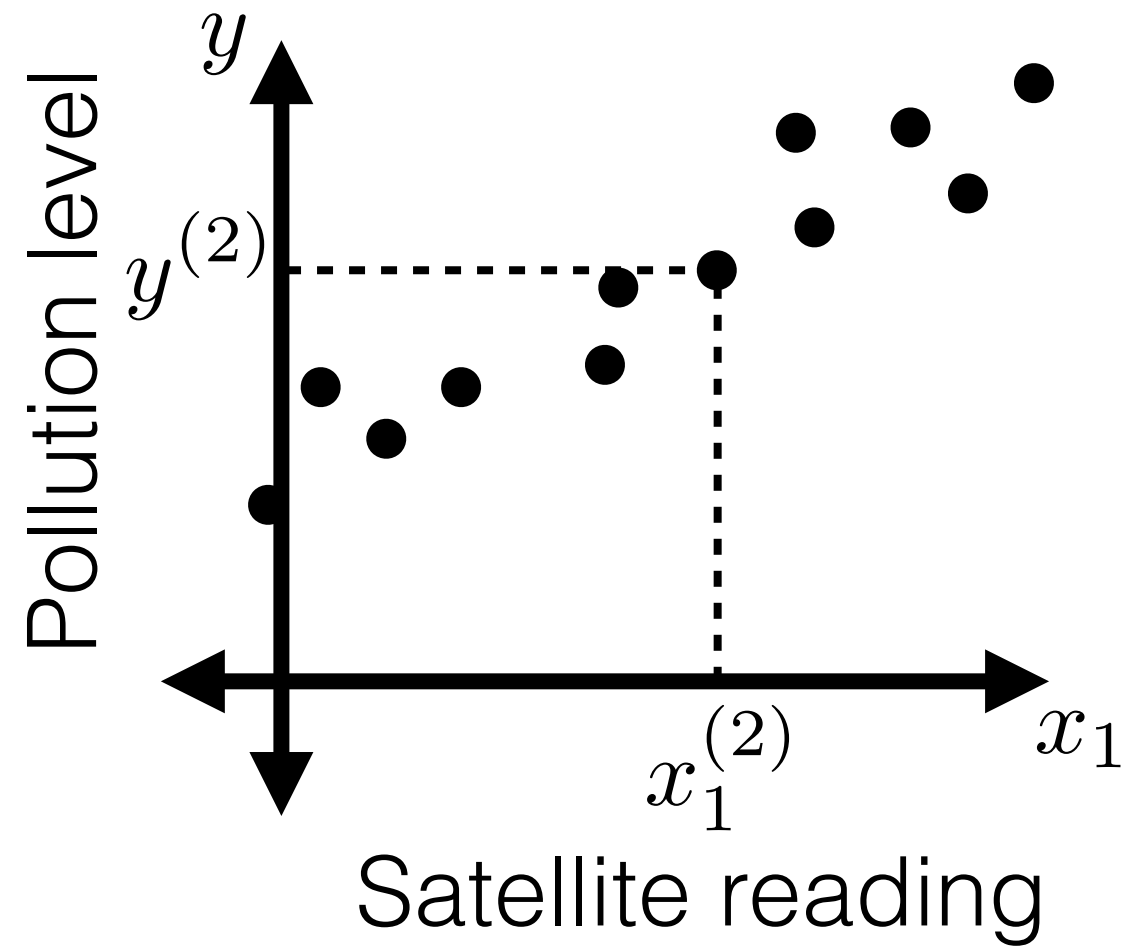
# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$



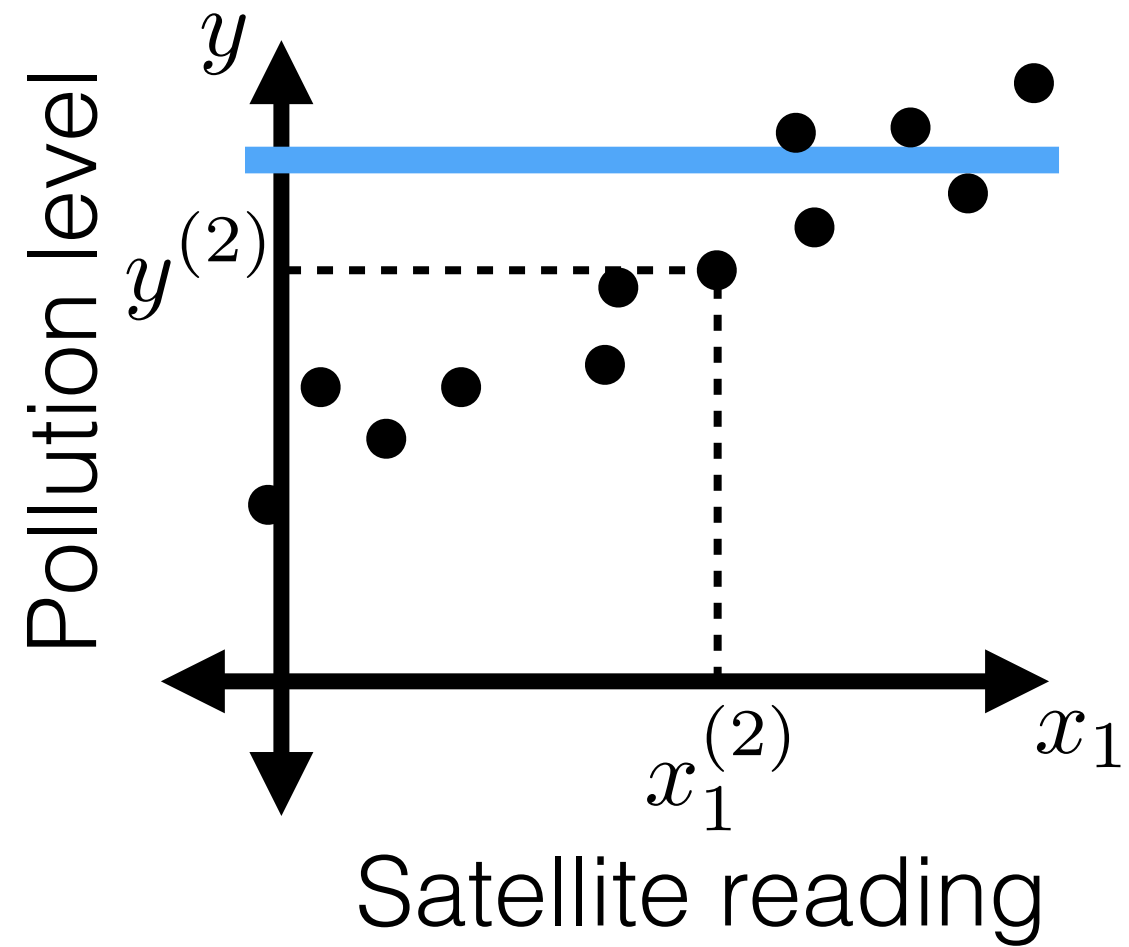
# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions



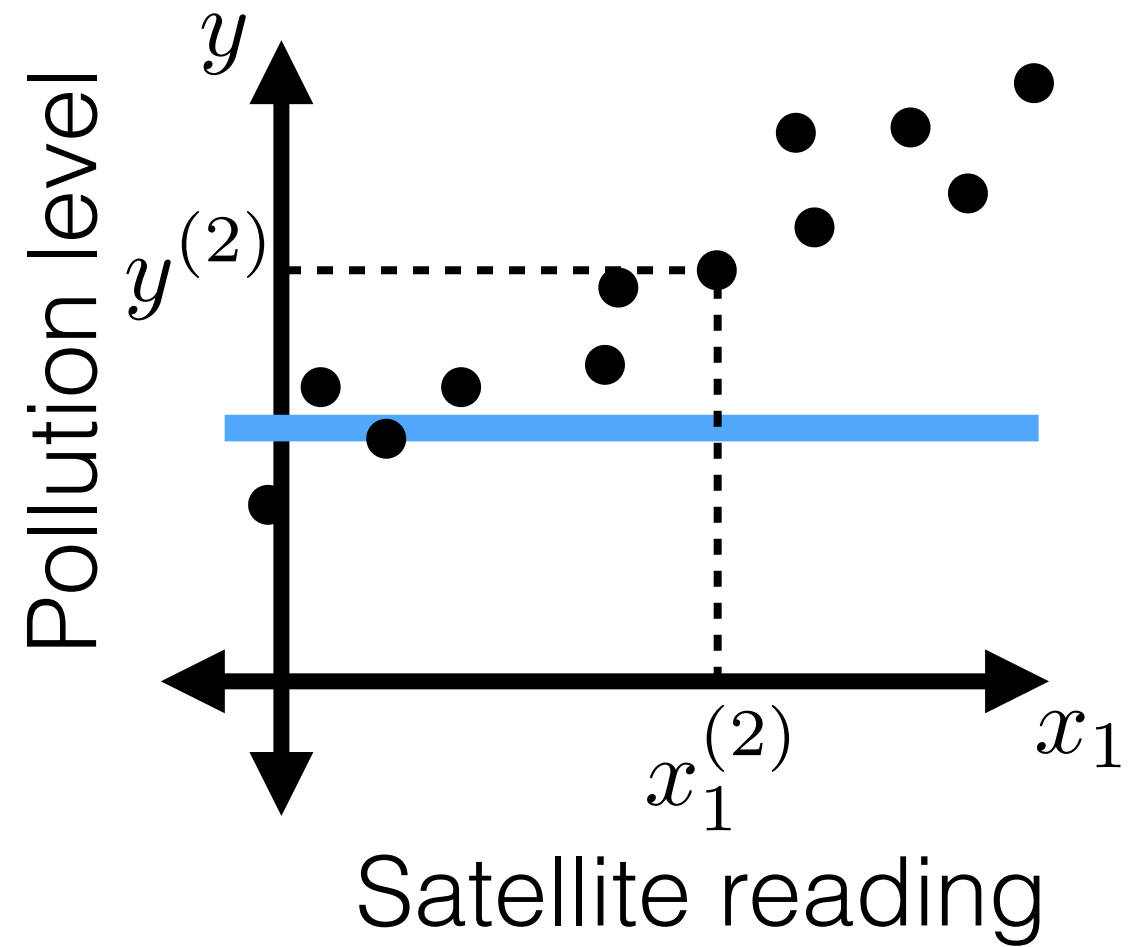
# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions



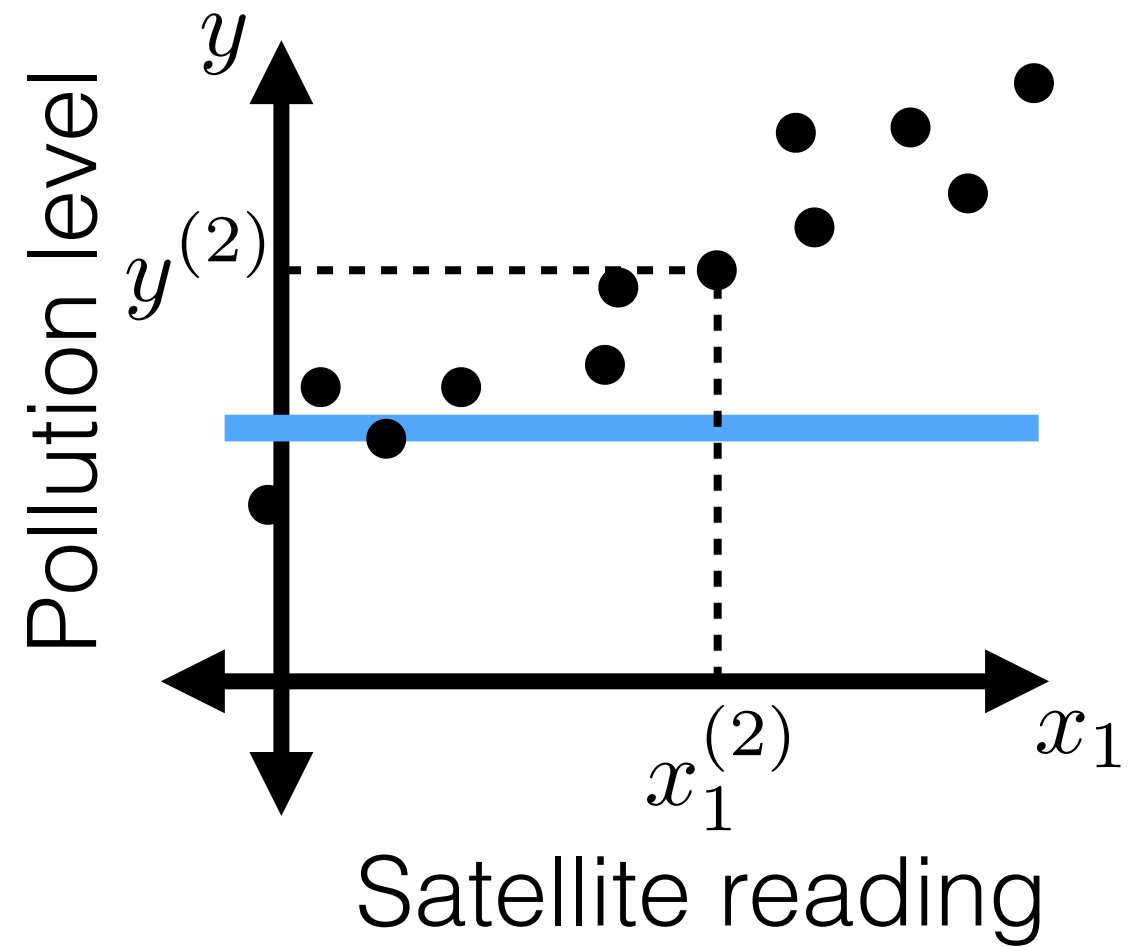
# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions

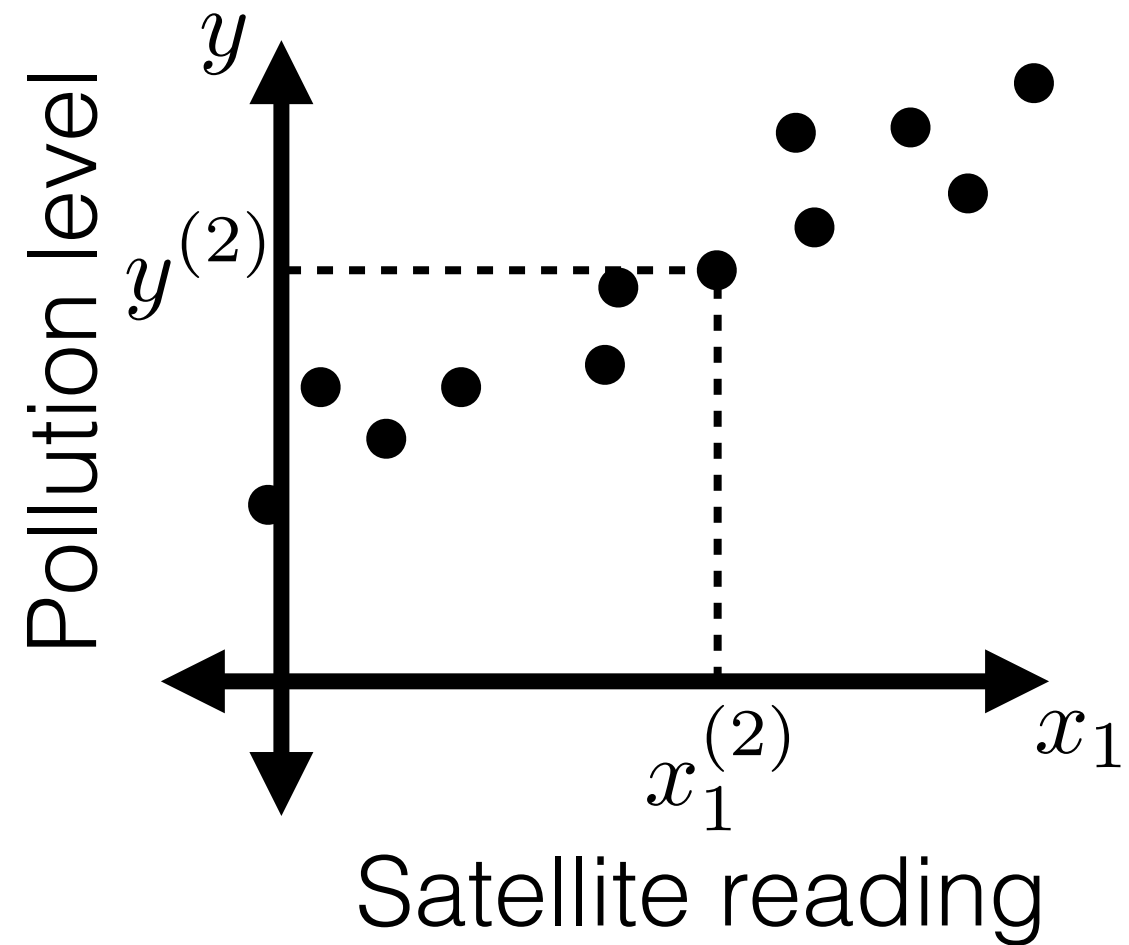




# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

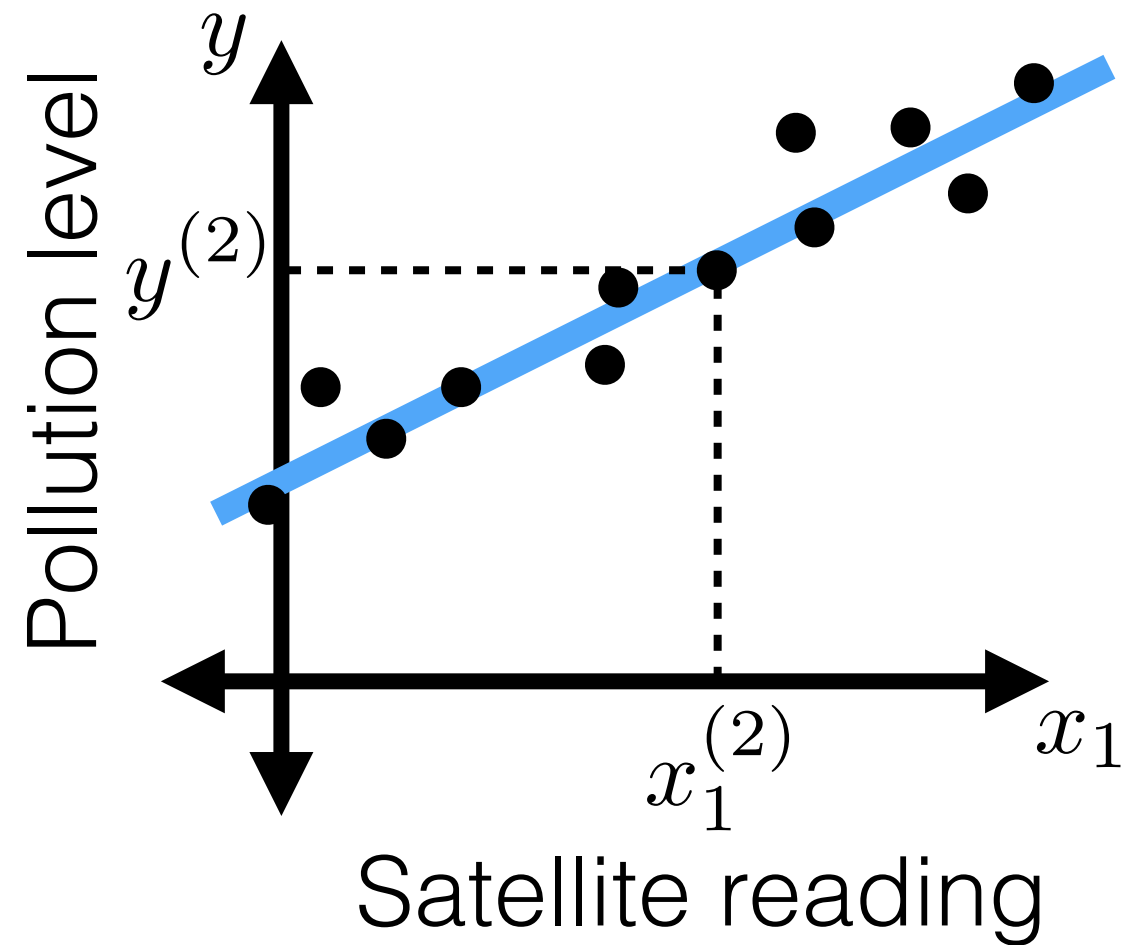
$$h(x) = \theta x + \theta_0$$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

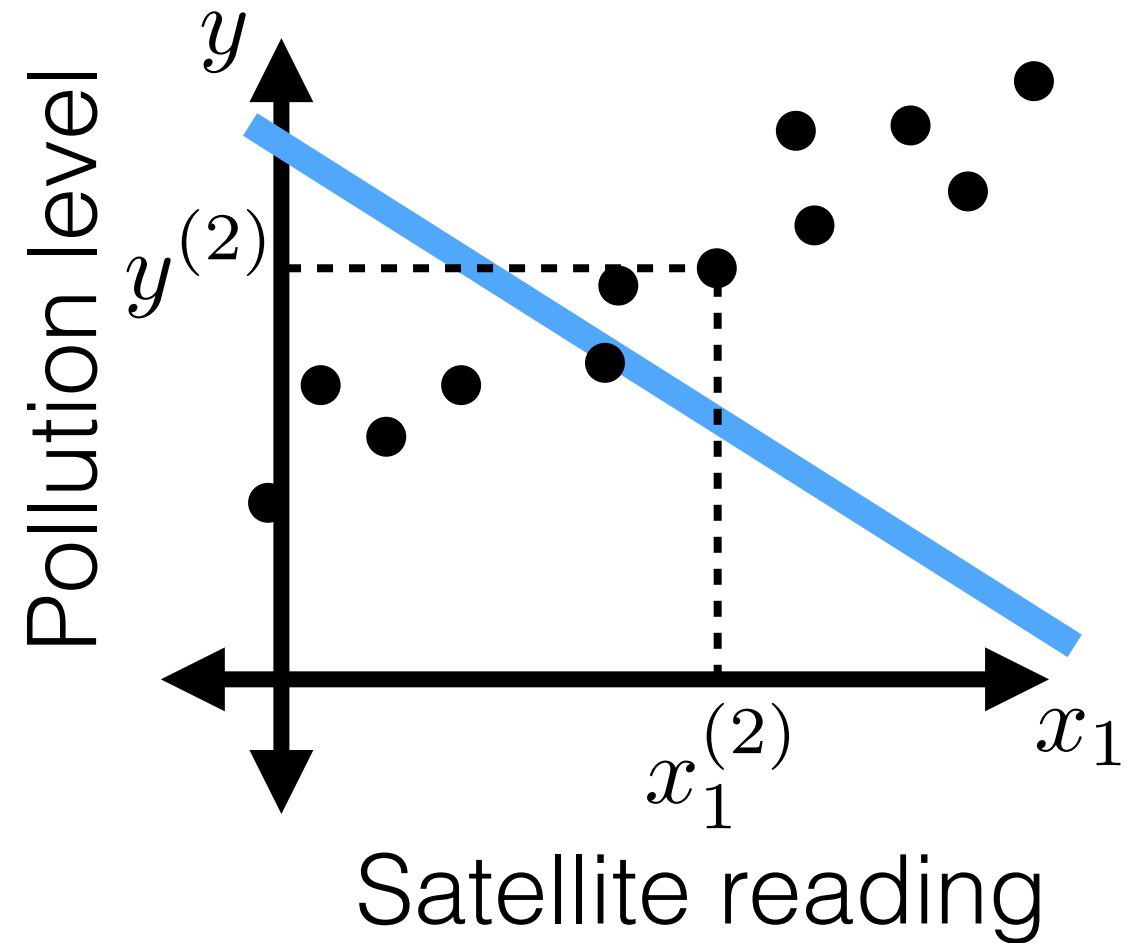
$$h(x) = \theta x + \theta_0$$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

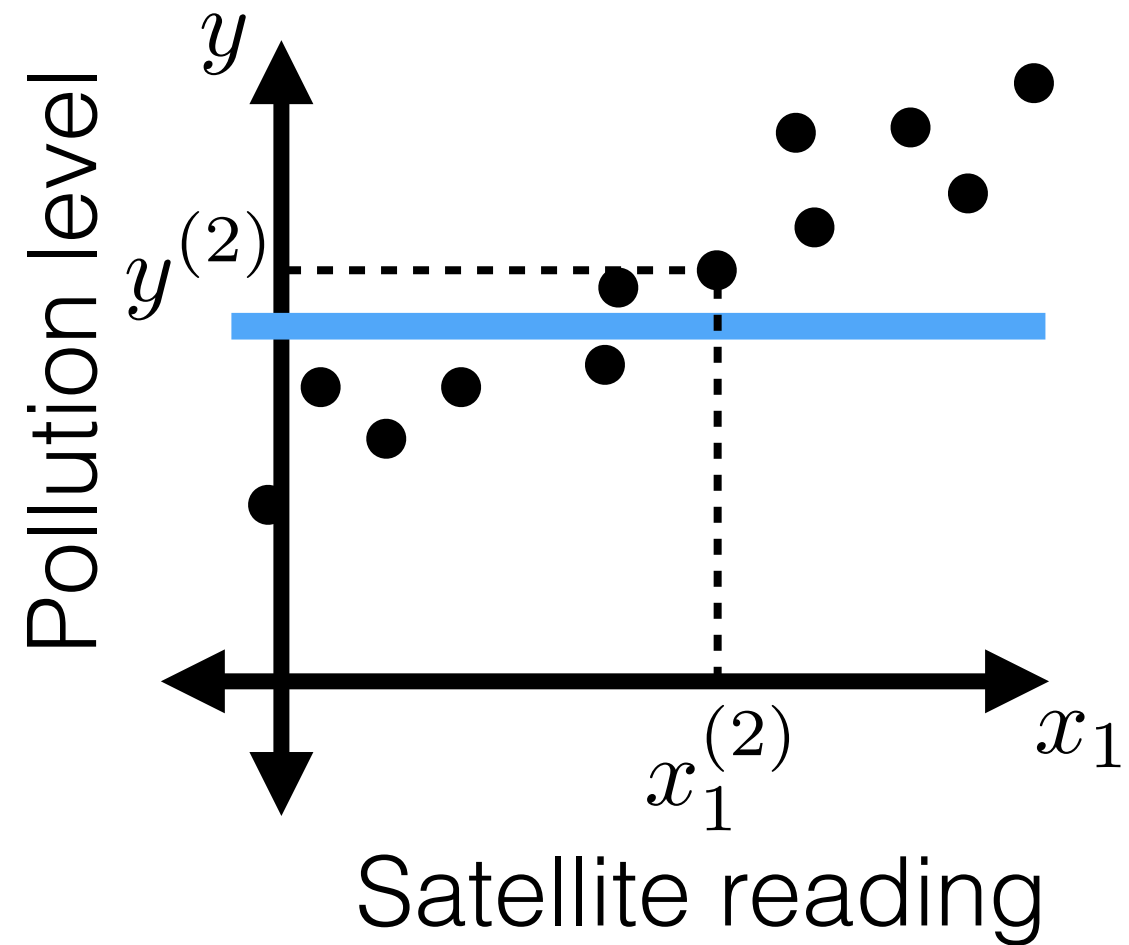
$$h(x) = \theta x + \theta_0$$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

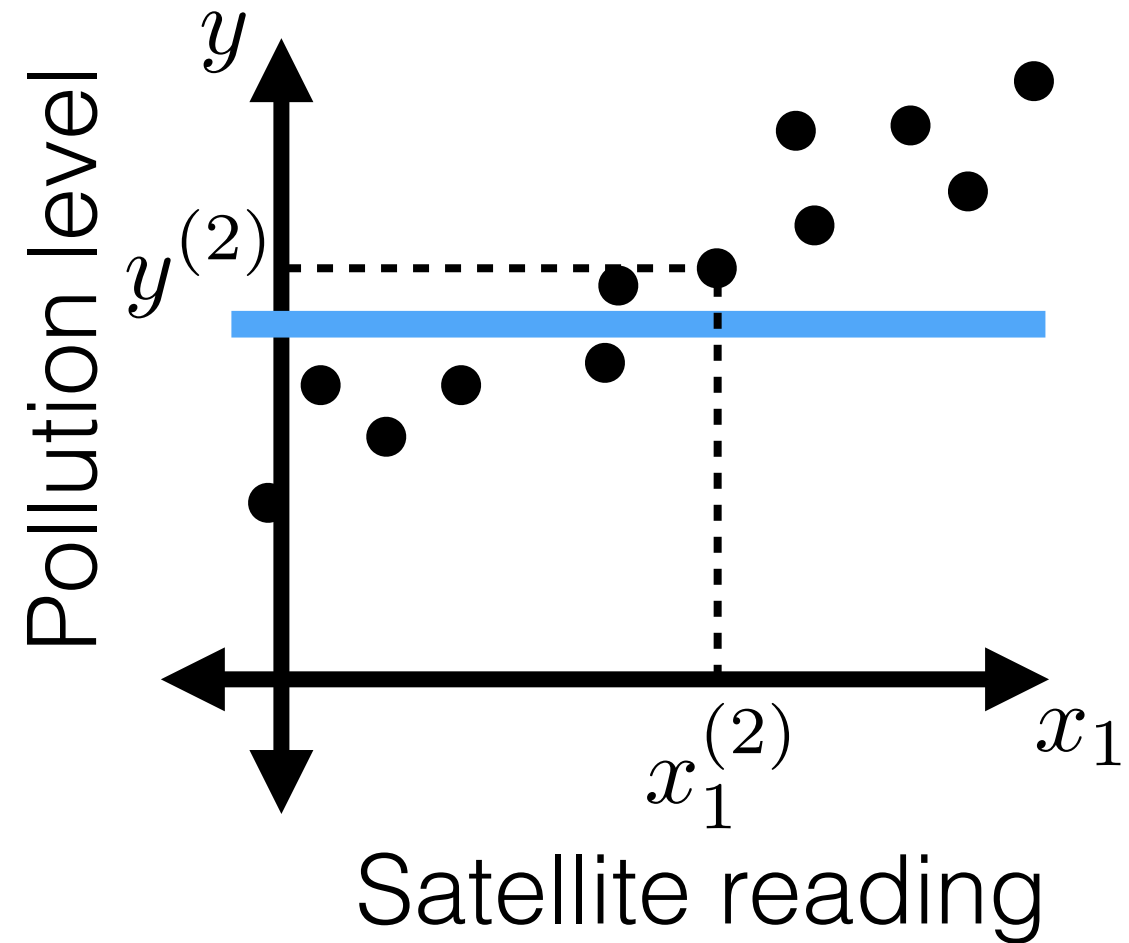
$$h(x) = \theta x + \theta_0$$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

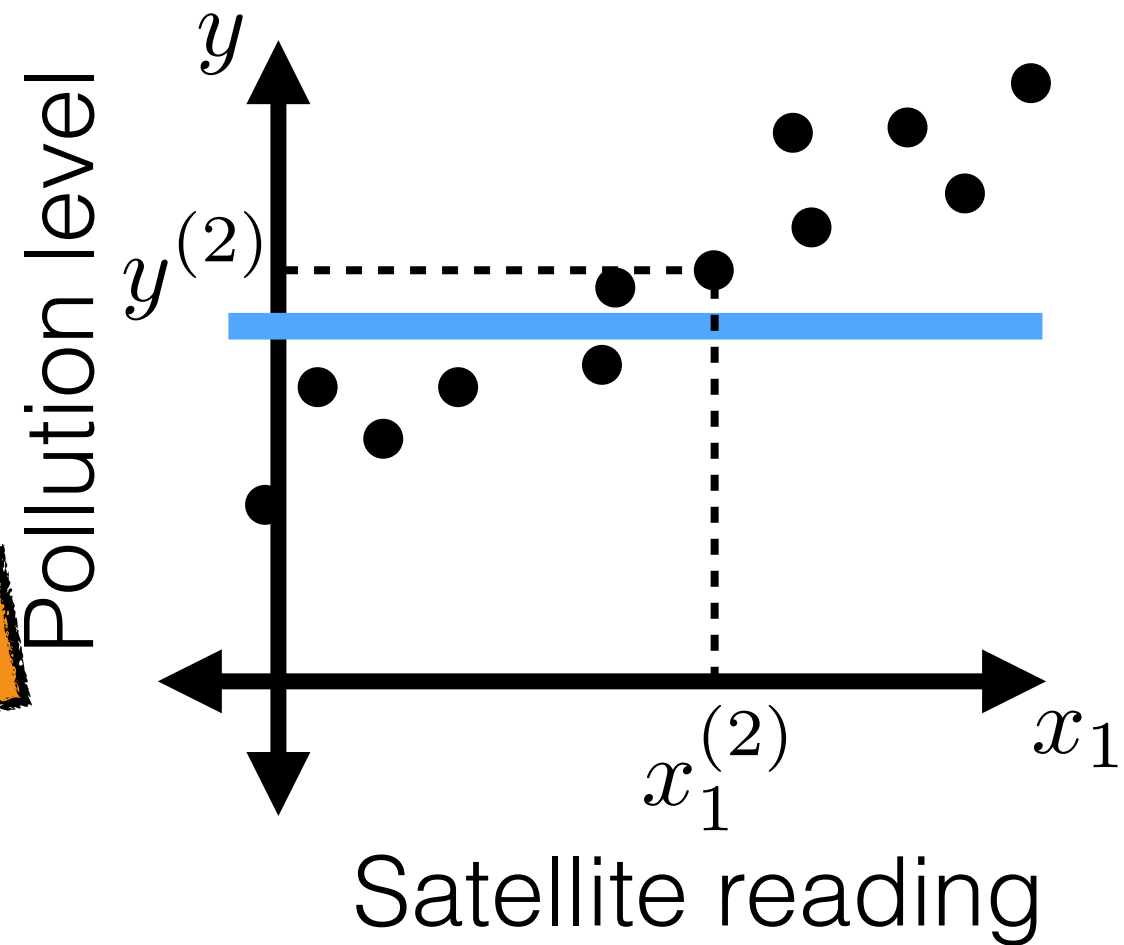


# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters



# Linear regressors

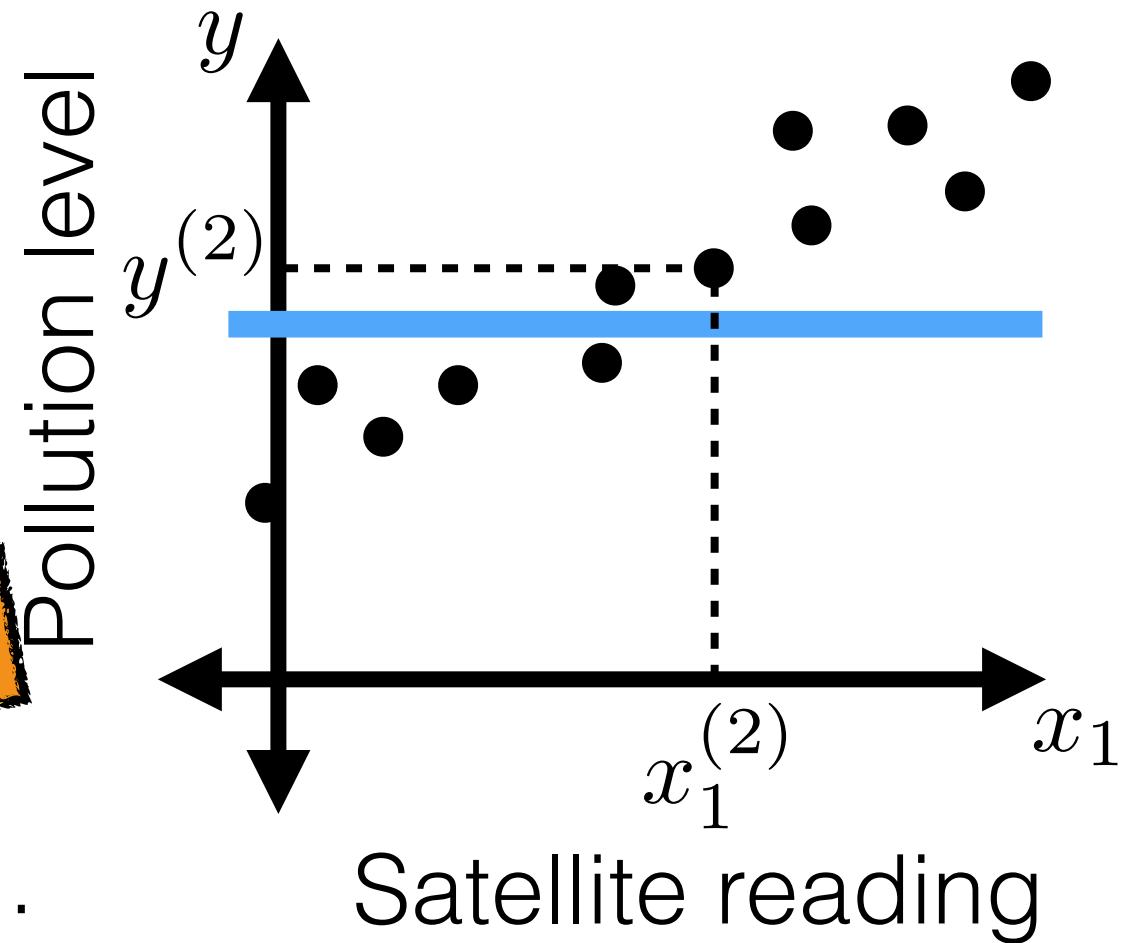
- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$h(x; \theta, \theta_0) = \theta_1 x_1 + \dots + \theta_d x_d + \theta_0$$



# Linear regressors

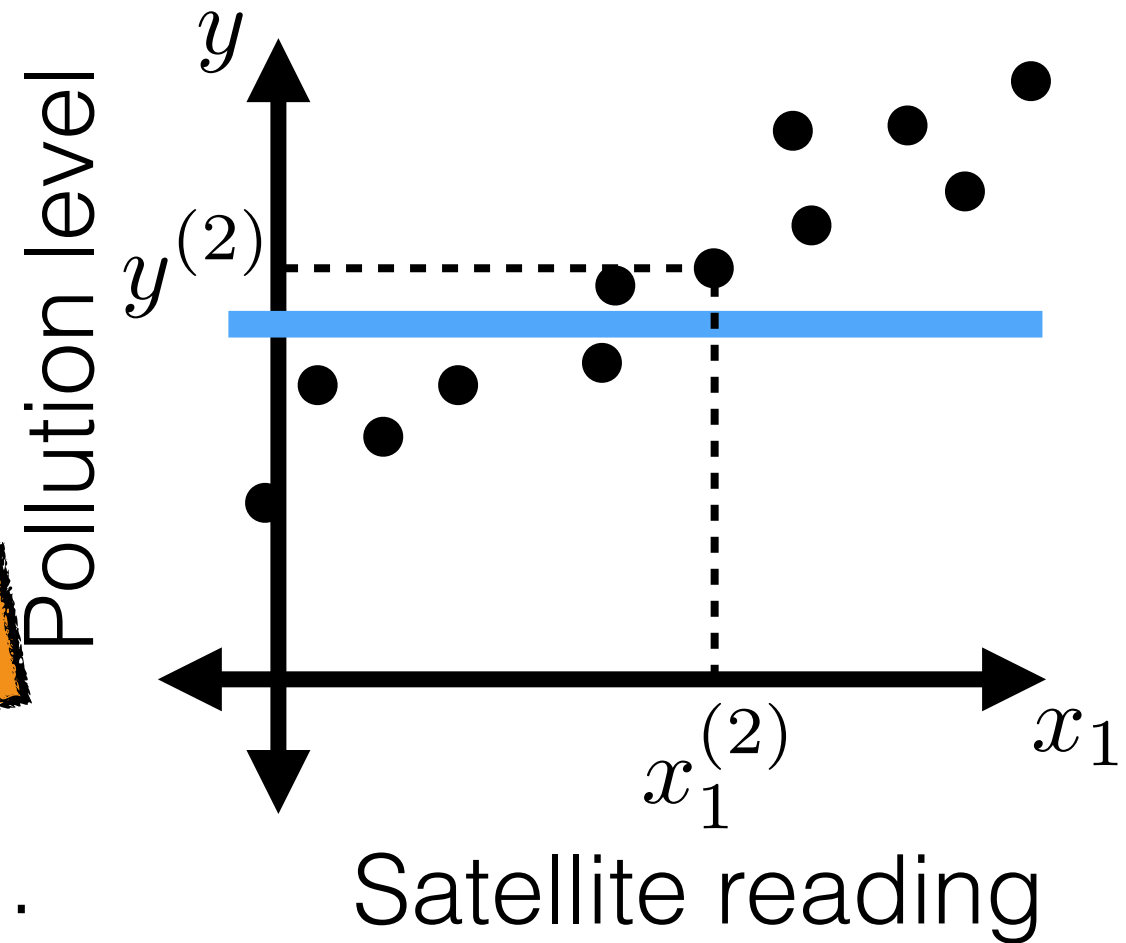
- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$





# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

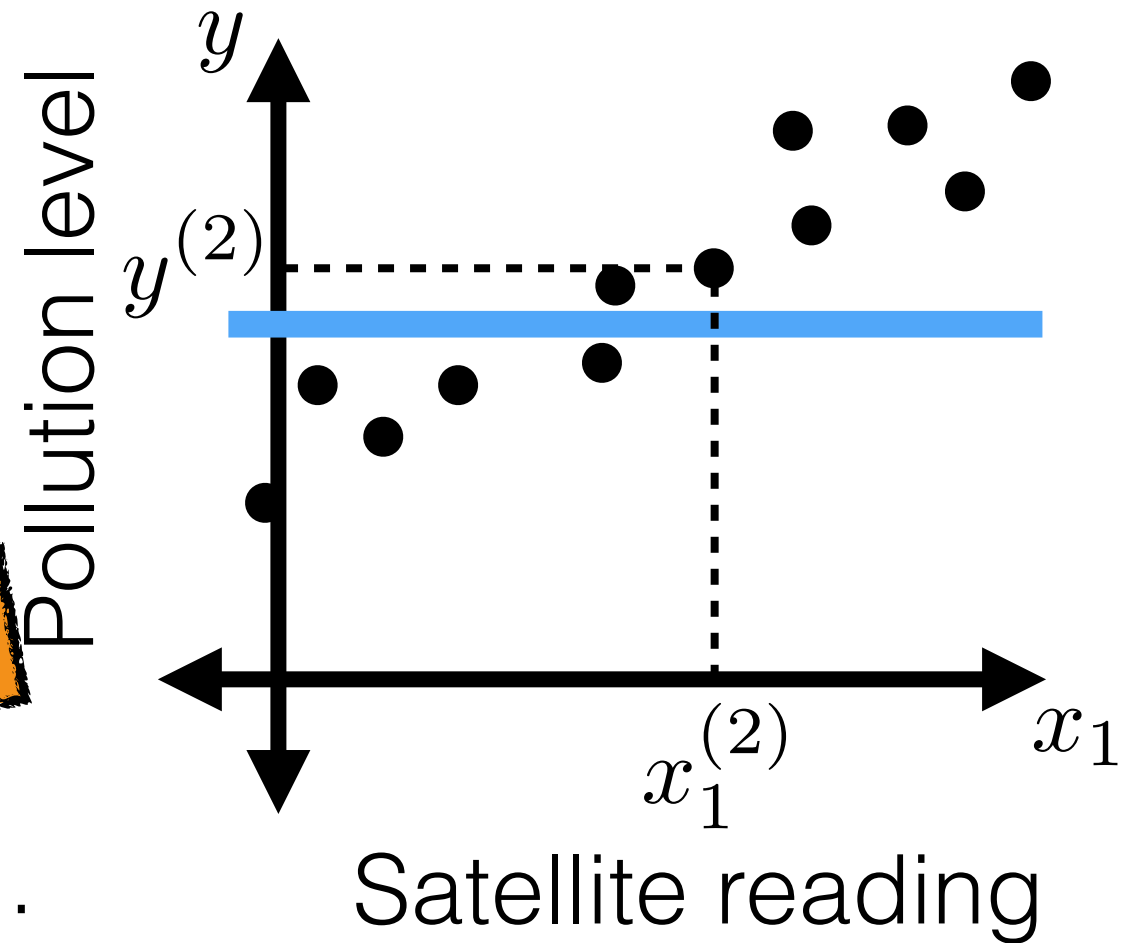
$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

$1 \times d, d \times 1$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

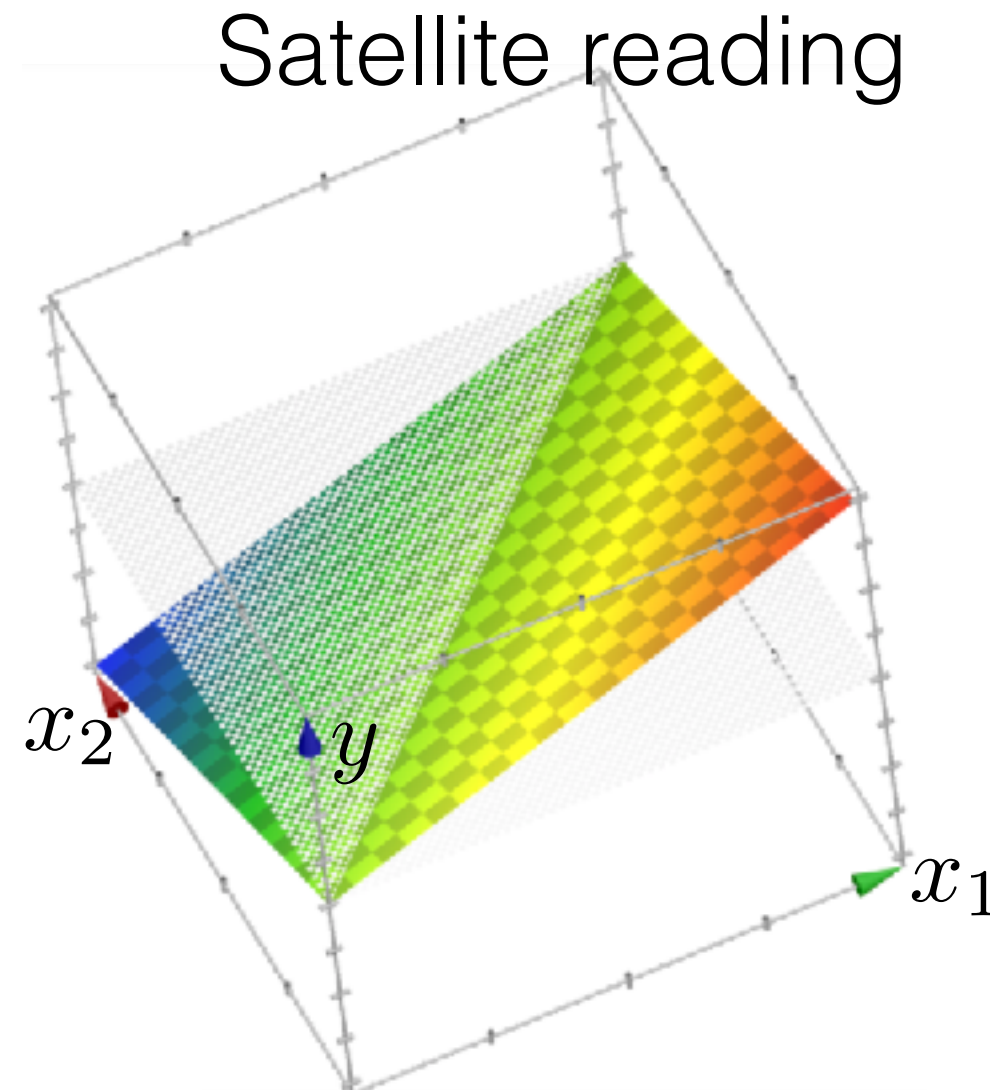
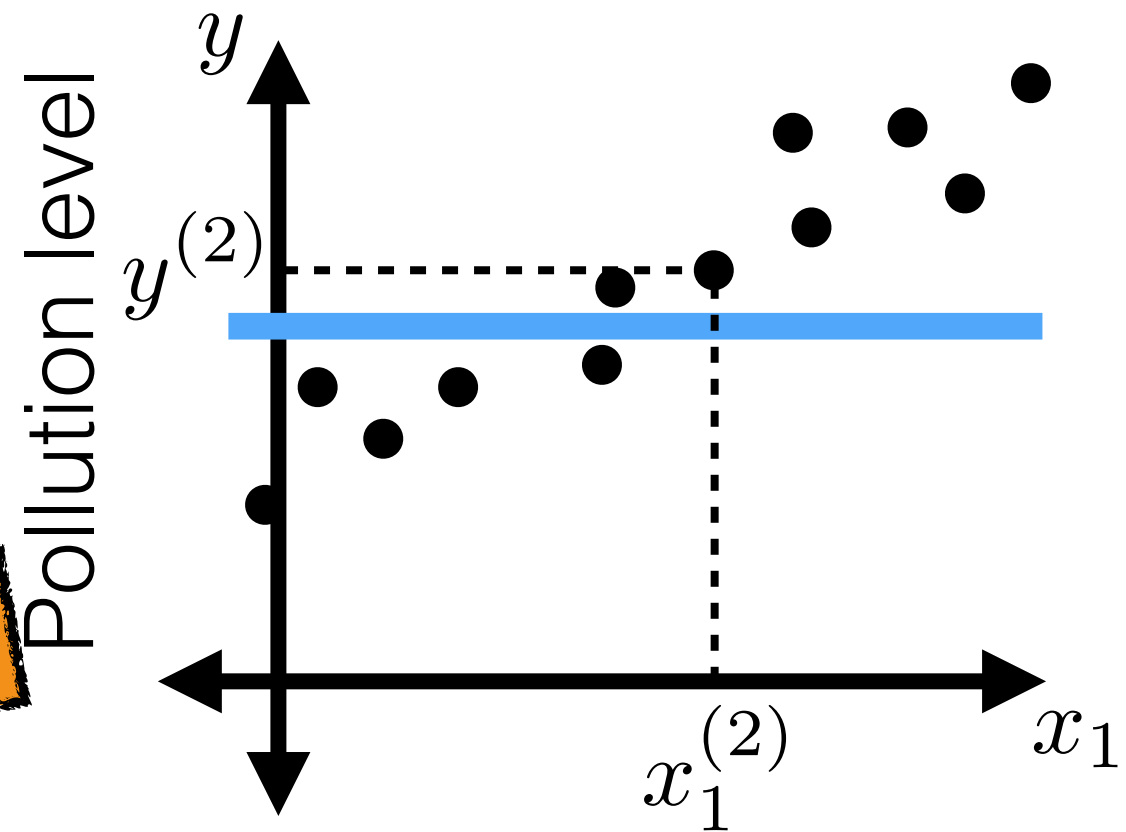
$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

$1 \times d, d \times 1$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

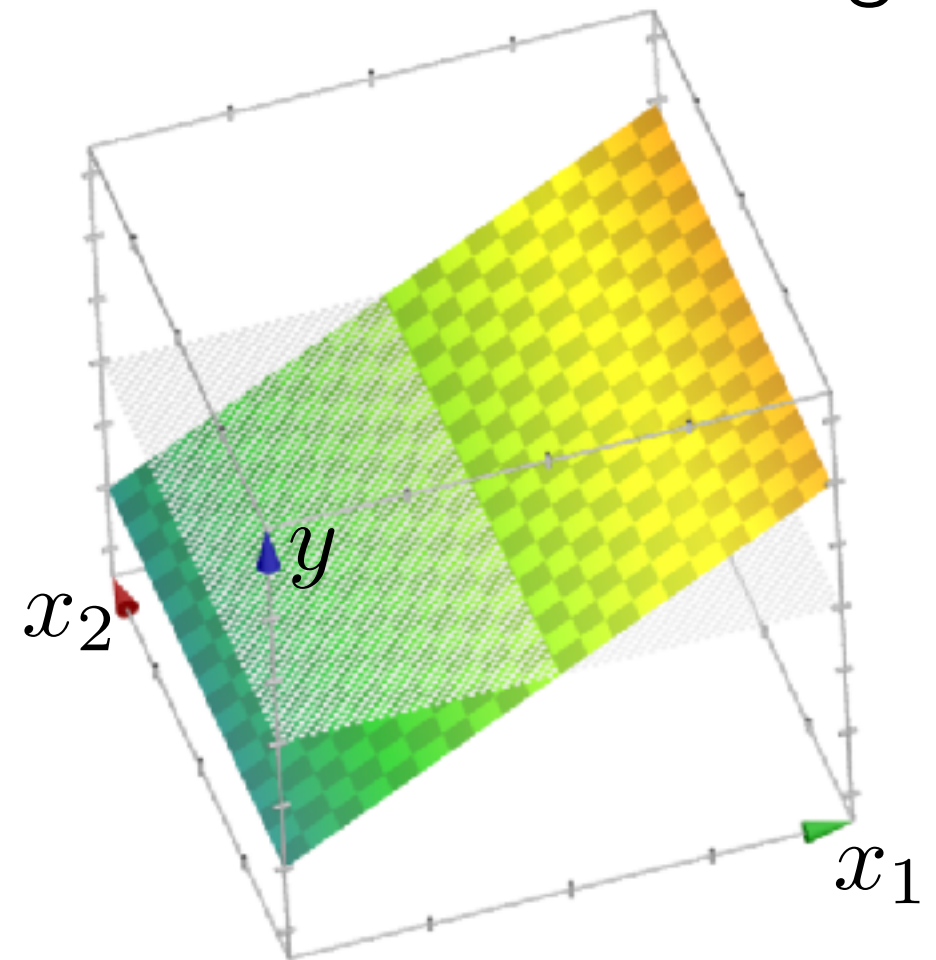
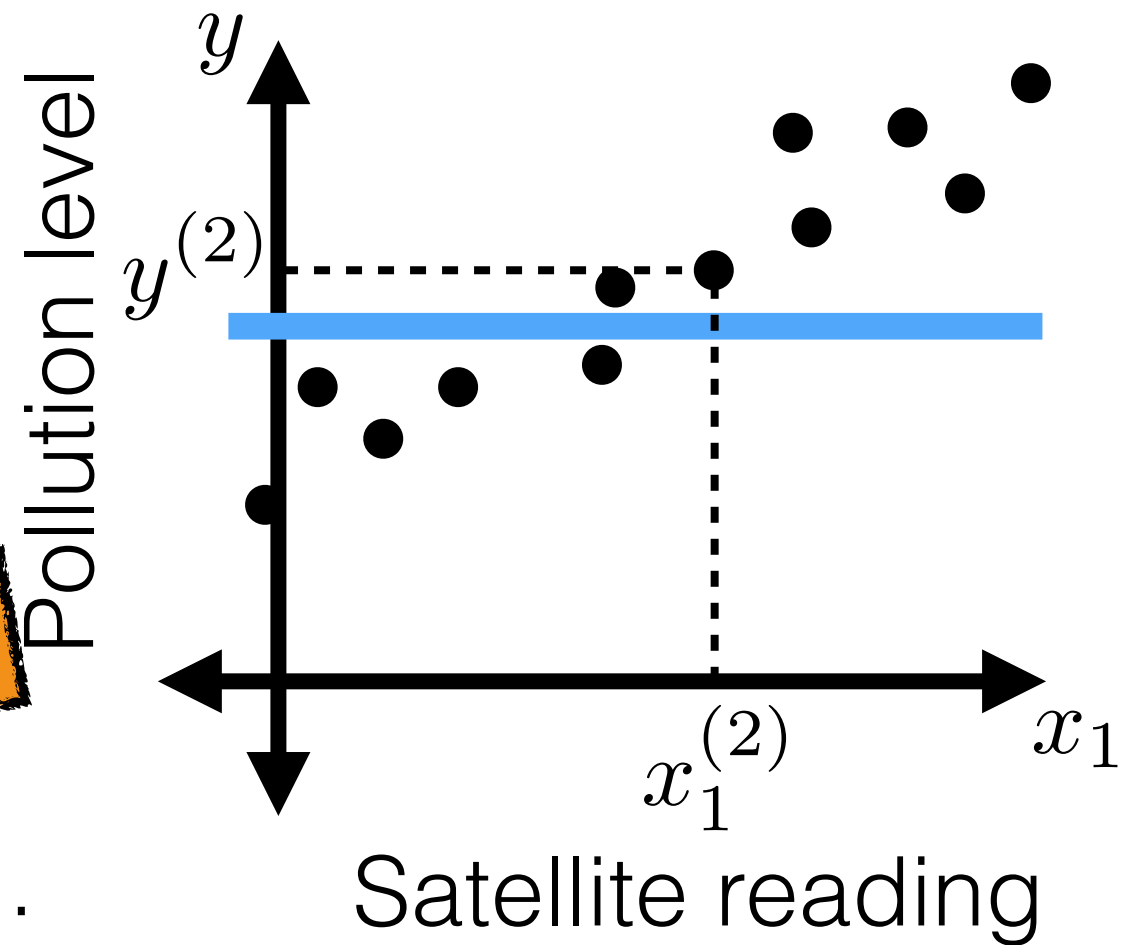
$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

$1 \times d, d \times 1$



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

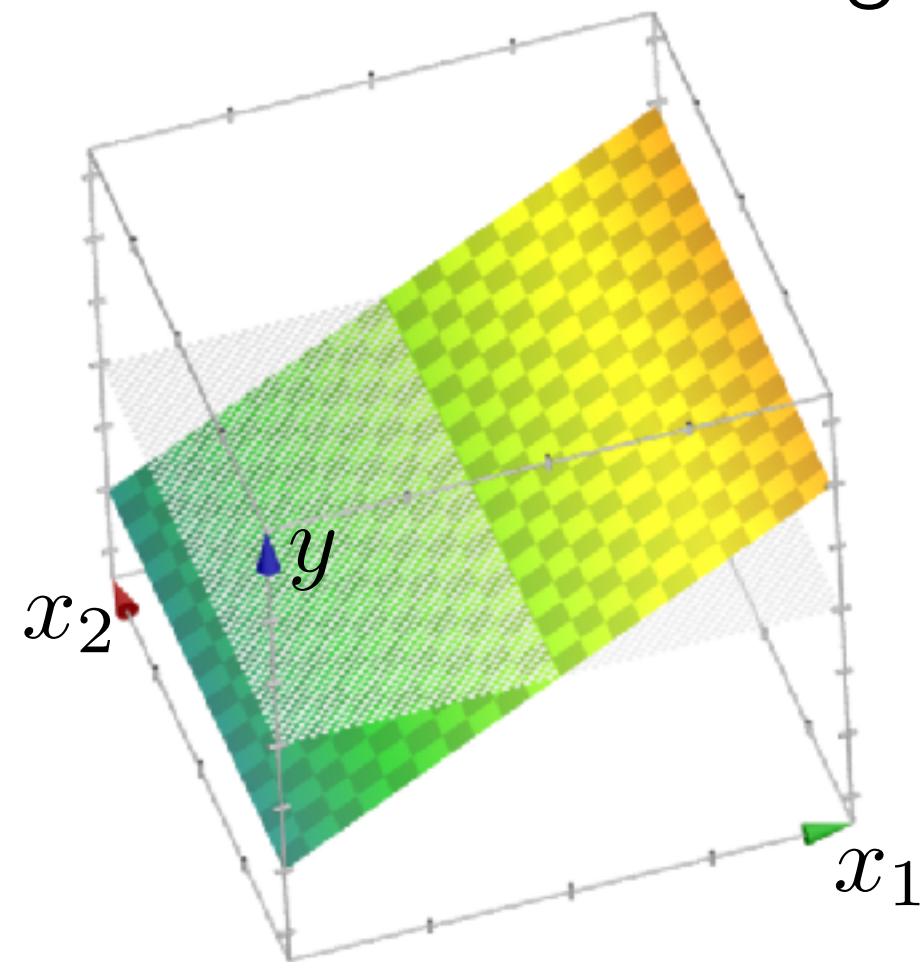
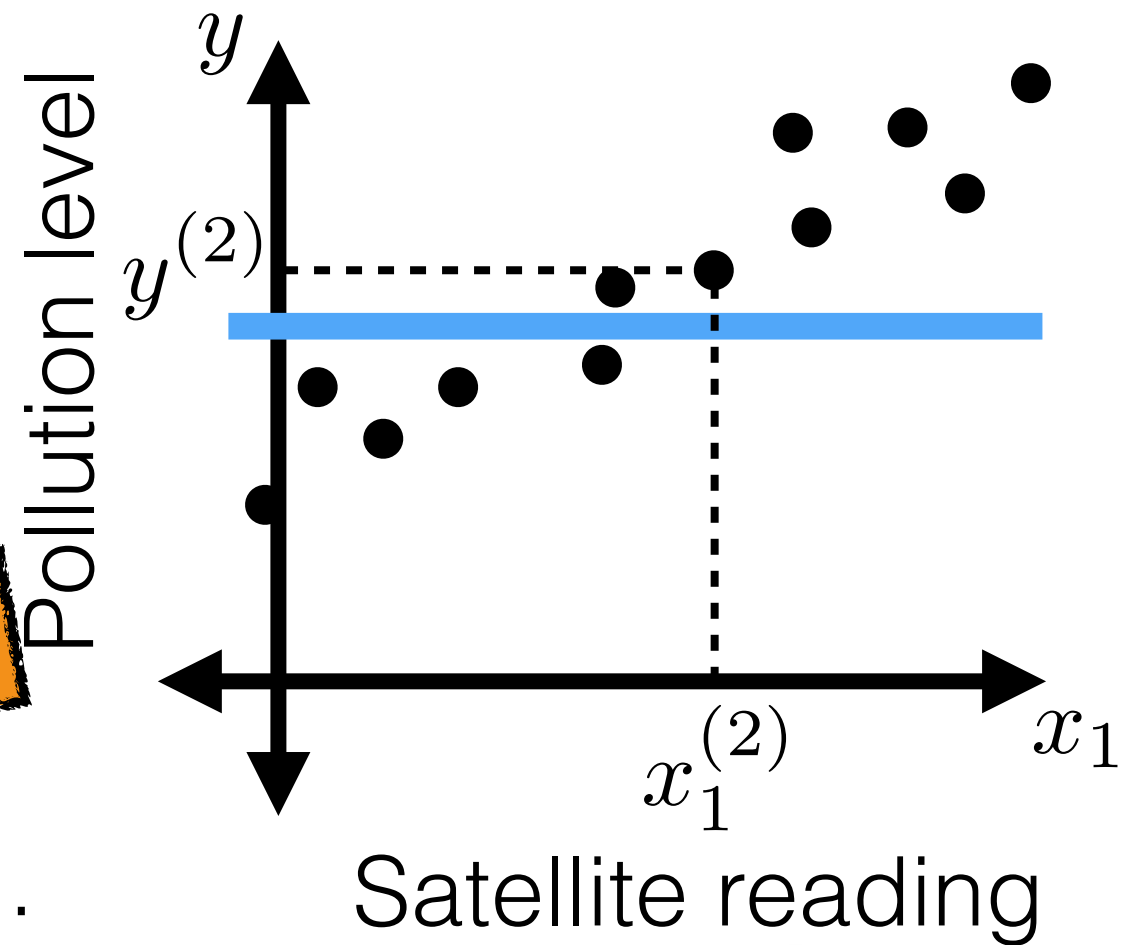
parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

$1 \times d, d \times 1$

OR



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

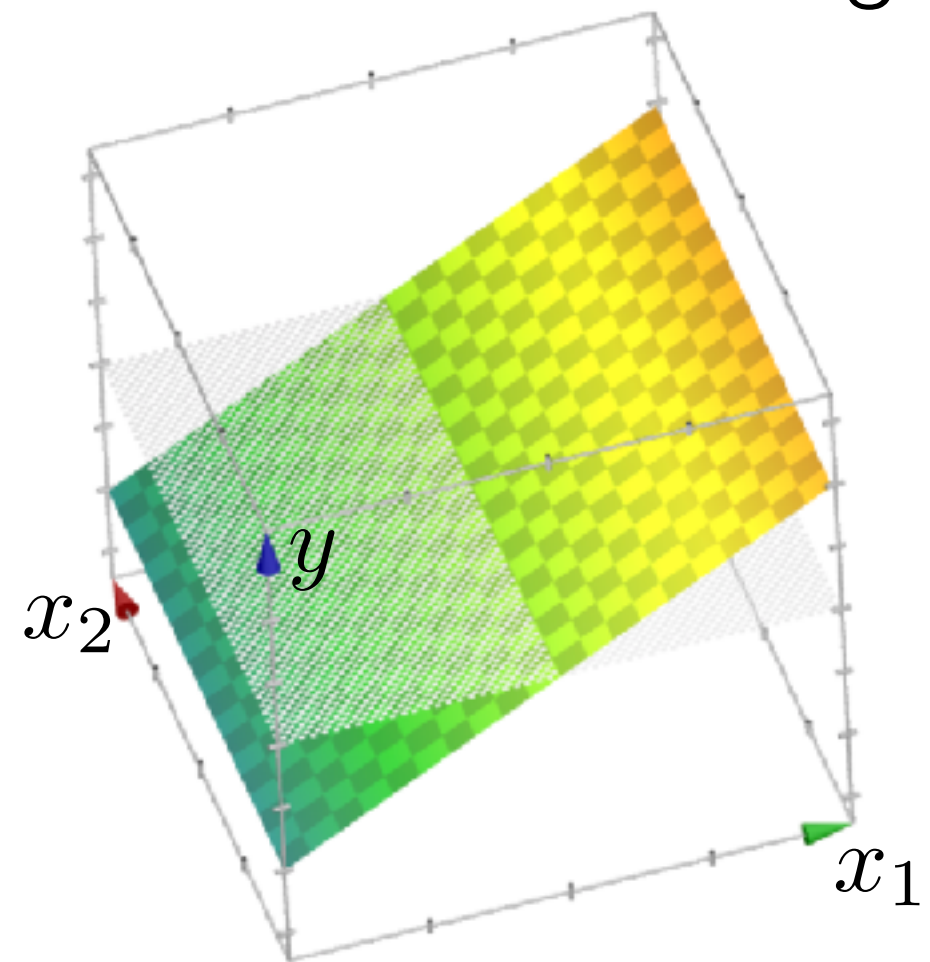
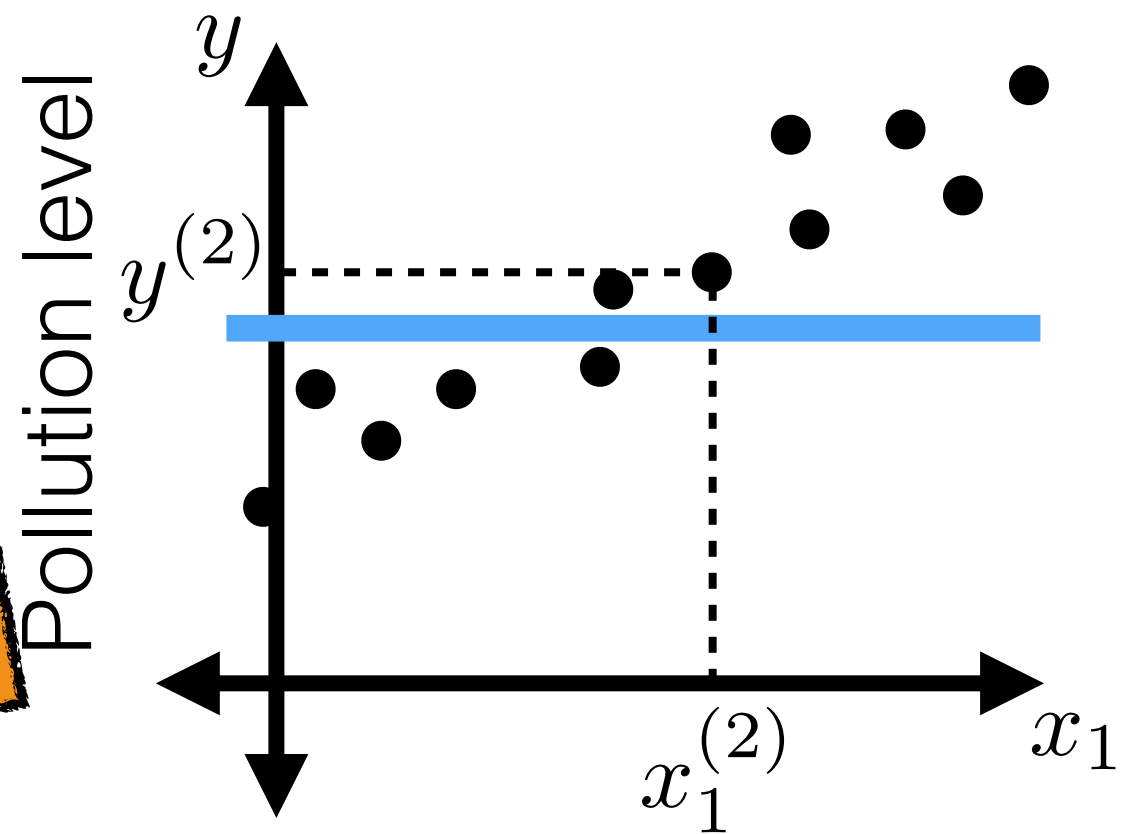
- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

1xd, dx1

OR

$$h(x) = \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1)$$





# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

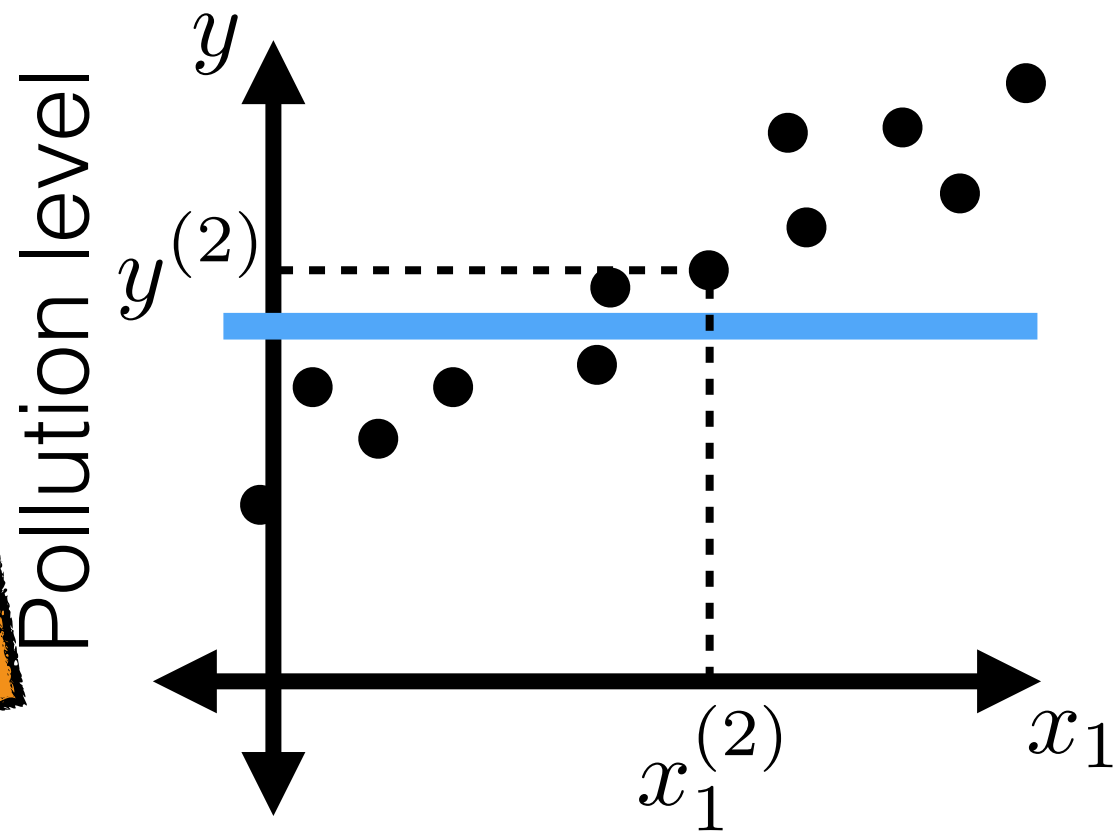
- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

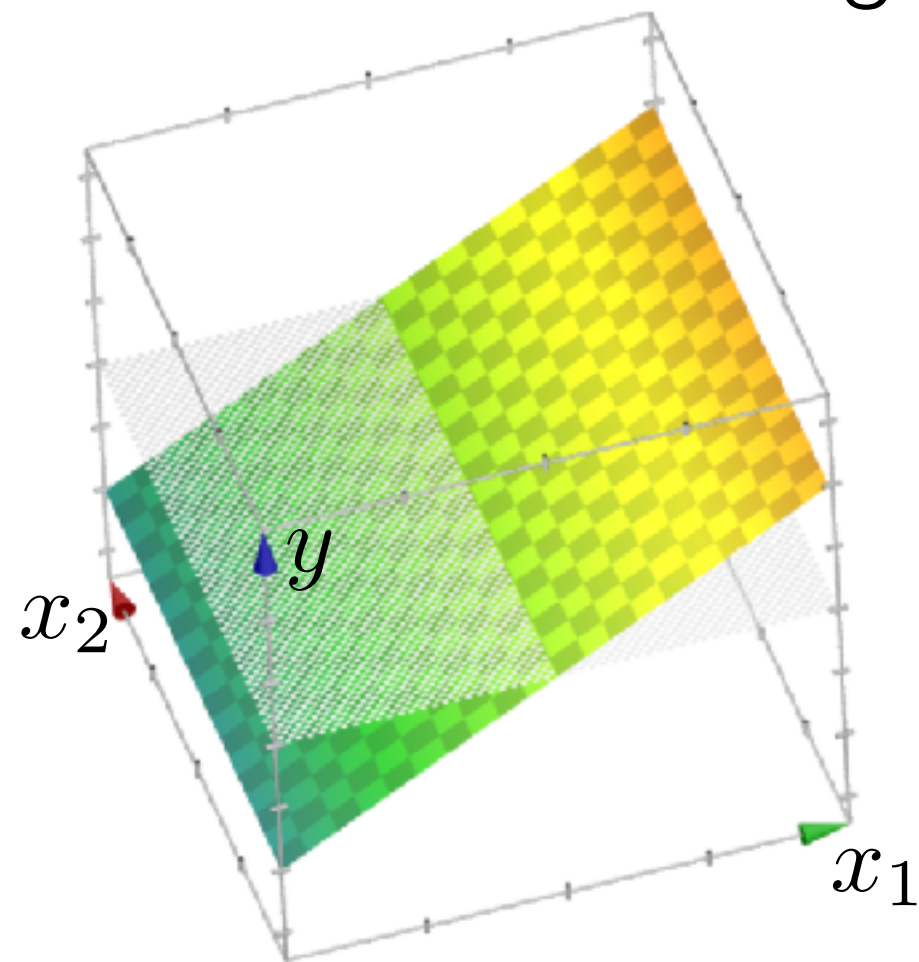
1xd, dx1

OR

$$\begin{aligned} h(x) &= \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1) \\ &= \theta^\top x \end{aligned}$$



Satellite reading



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$h(x; \theta, \theta_0) = \theta_1 x_1 + \dots + \theta_d x_d + \theta_0$$

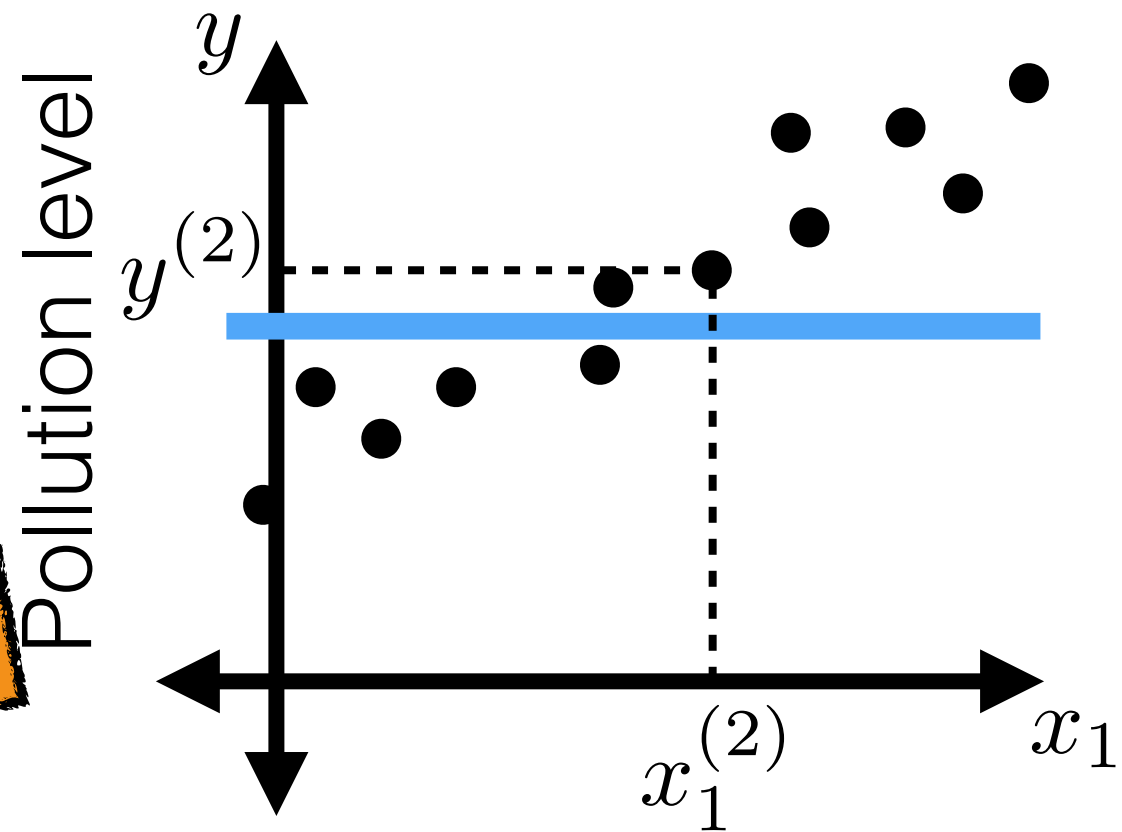
$$= \theta^\top x + \theta_0$$

$1 \times d, d \times 1$

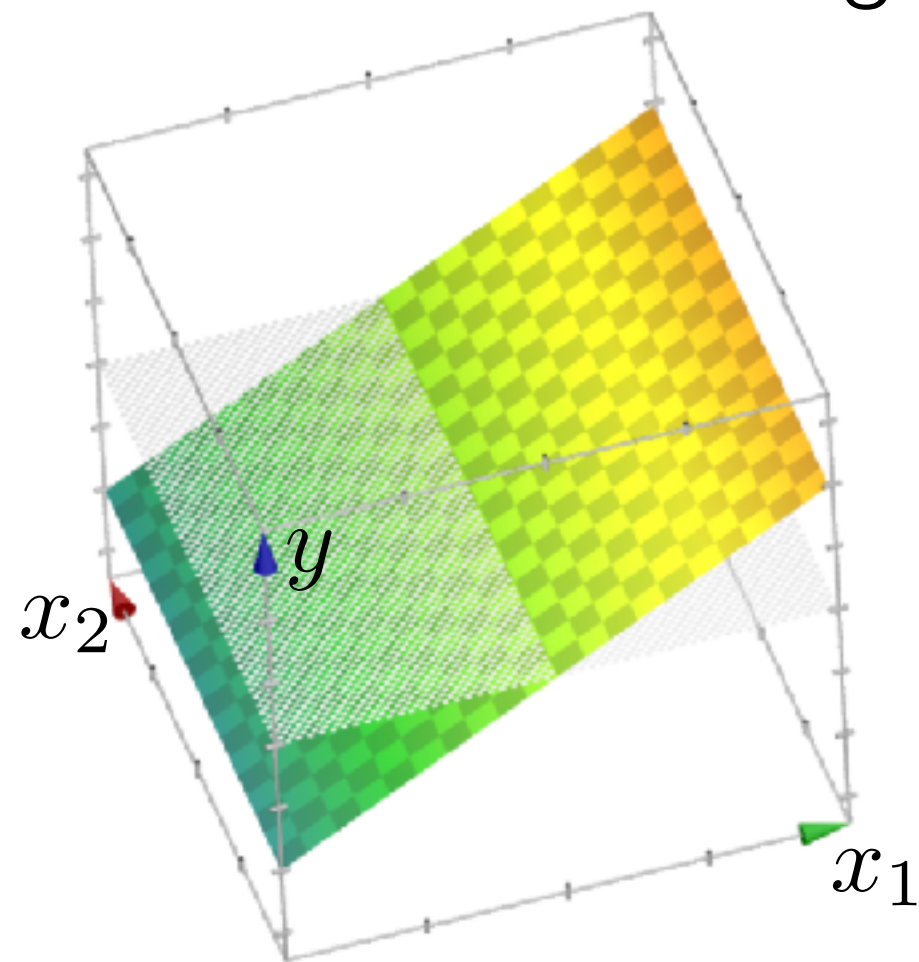
OR

$$h(x; \theta) = \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1)$$

$$= \theta^\top x$$



Satellite reading



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

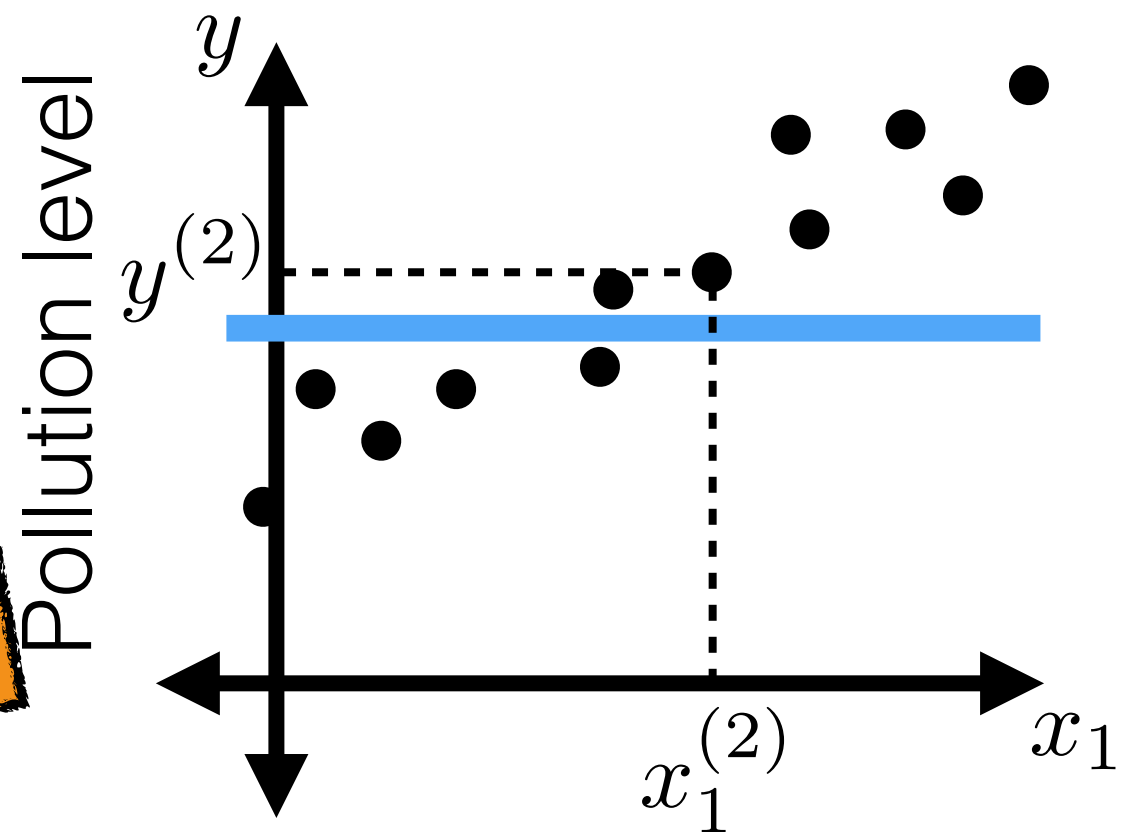
$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

1xd, dx1

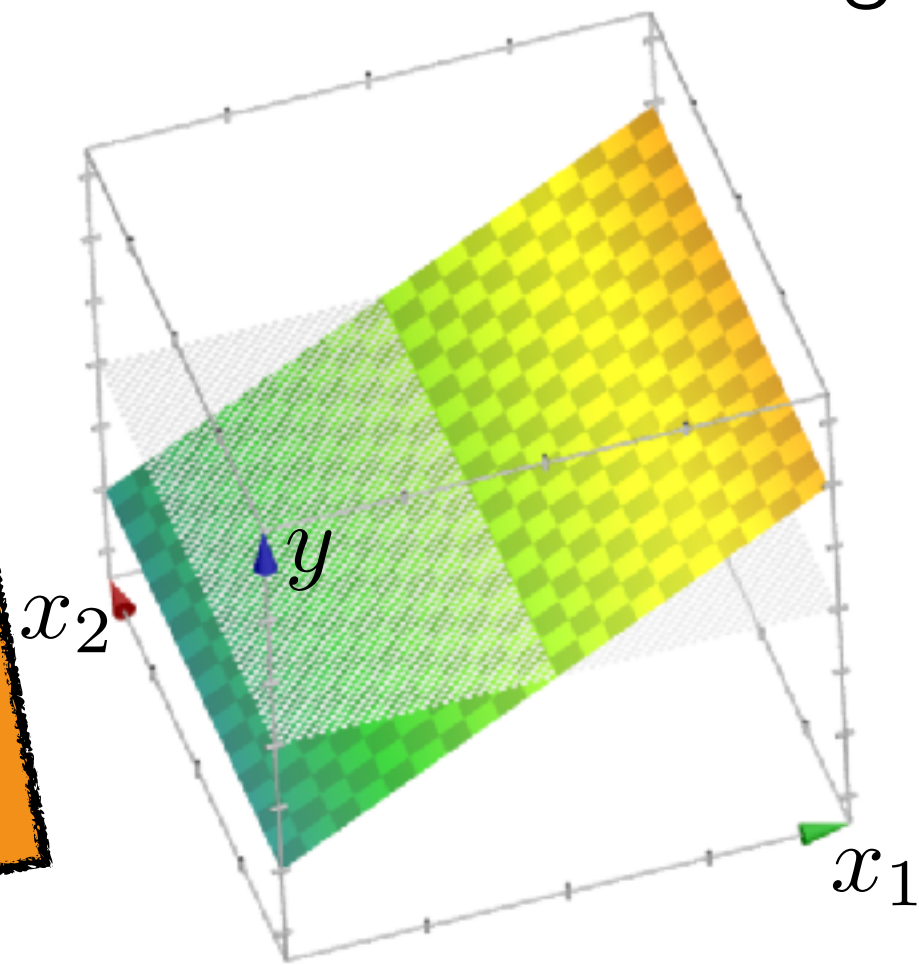
OR

$$\begin{aligned} h(x; \theta) &= \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1) \\ &= \theta^\top x \end{aligned}$$

Notational trick: not the same  $\theta$  &  $x$ !



Satellite reading





# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

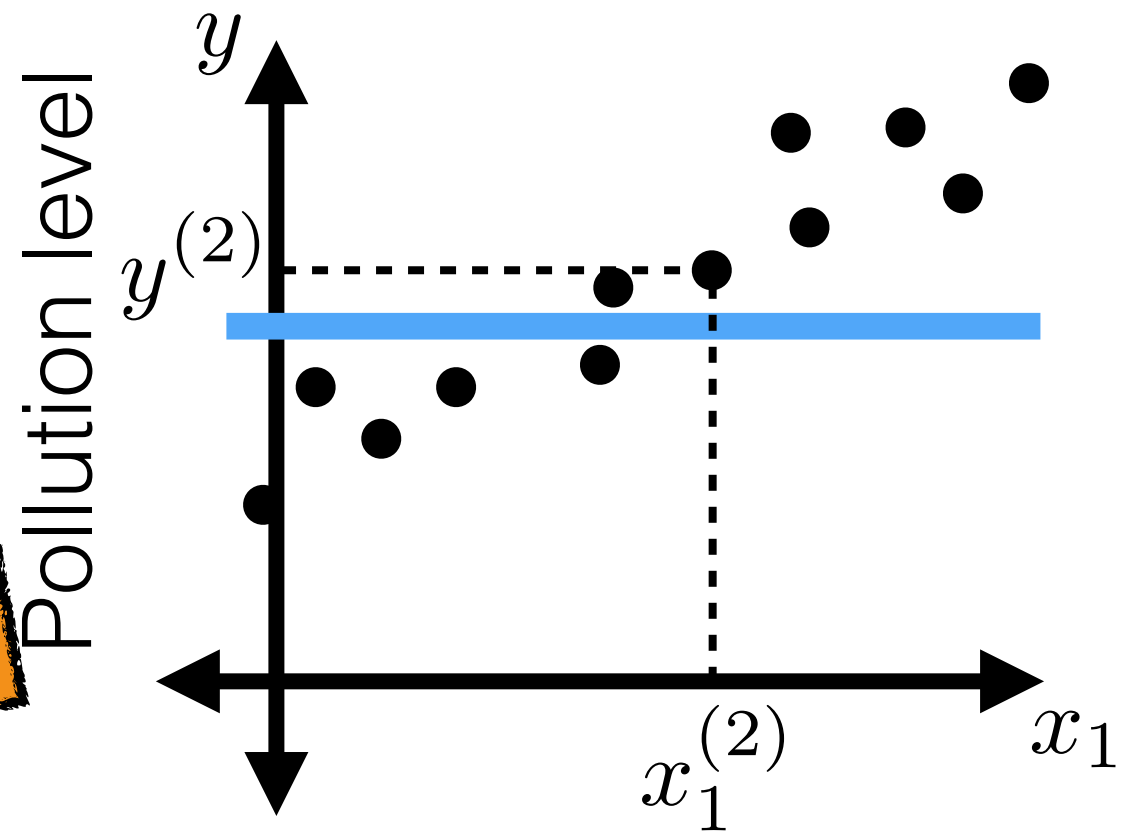
$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

1x2, 2x1

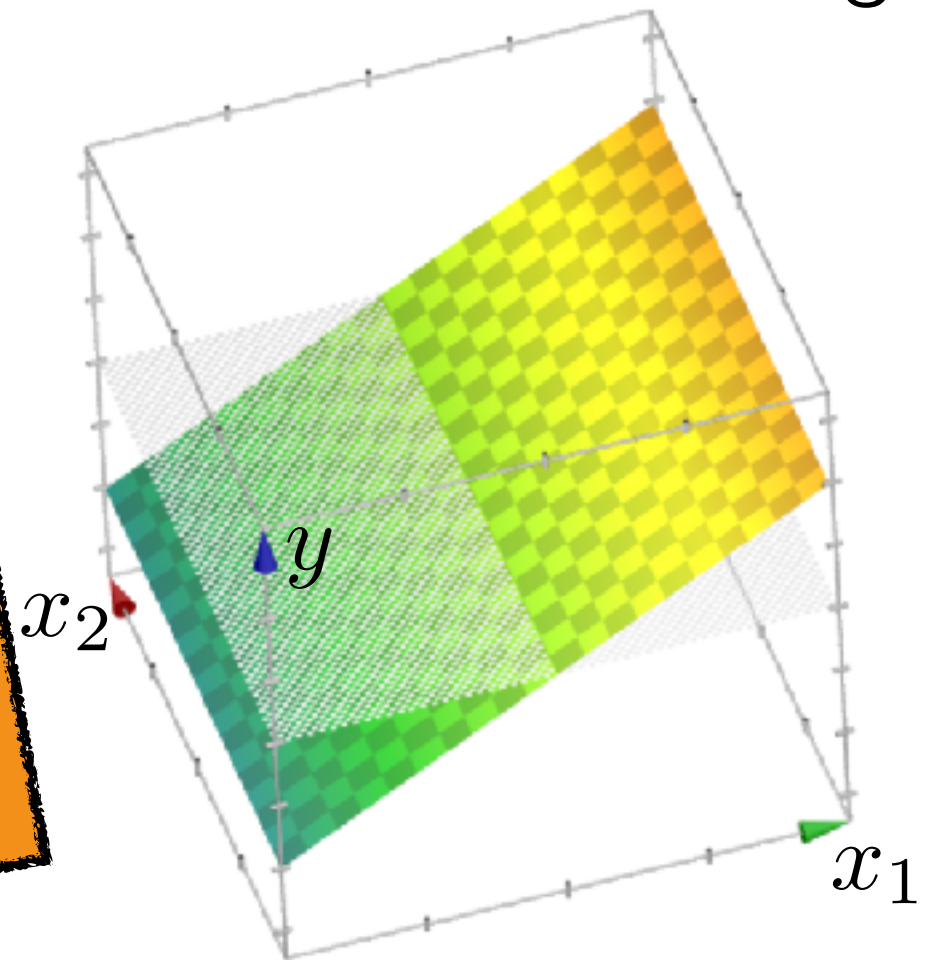
OR

$$\begin{aligned} h(x; \theta) &= \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1) \\ &= \theta^\top x \end{aligned}$$

Notational trick: not the same  $\theta$  &  $x$ !



Satellite reading



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

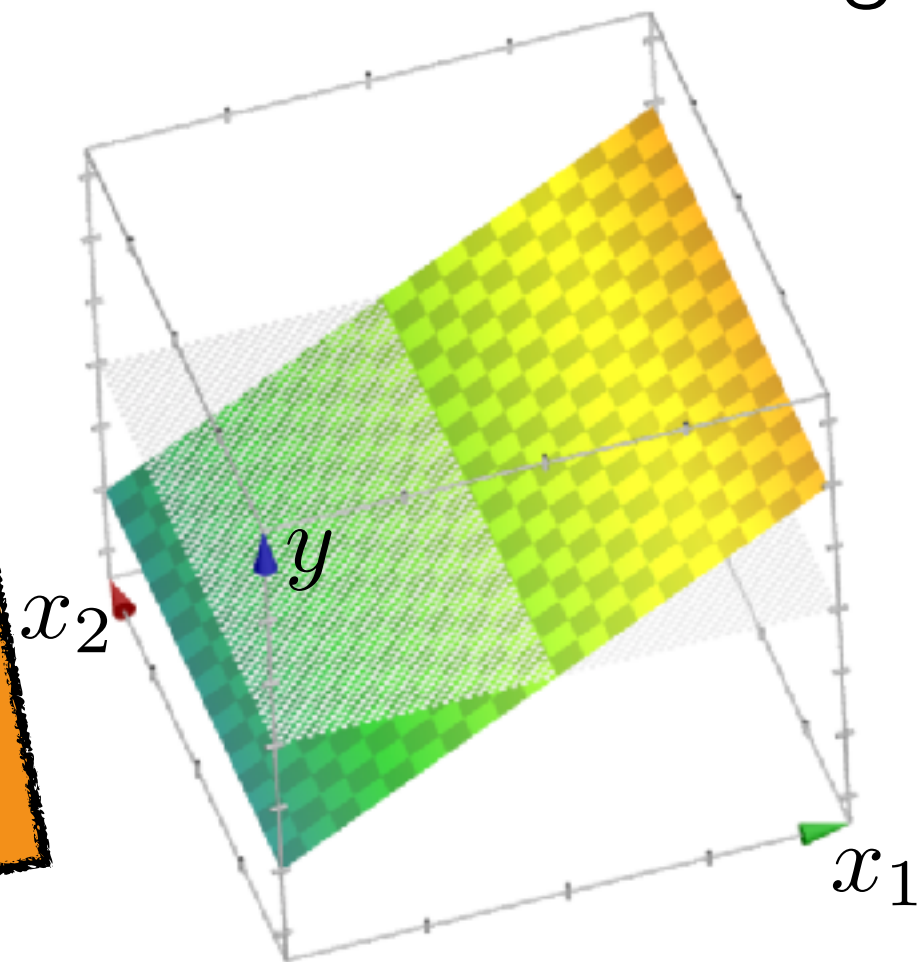
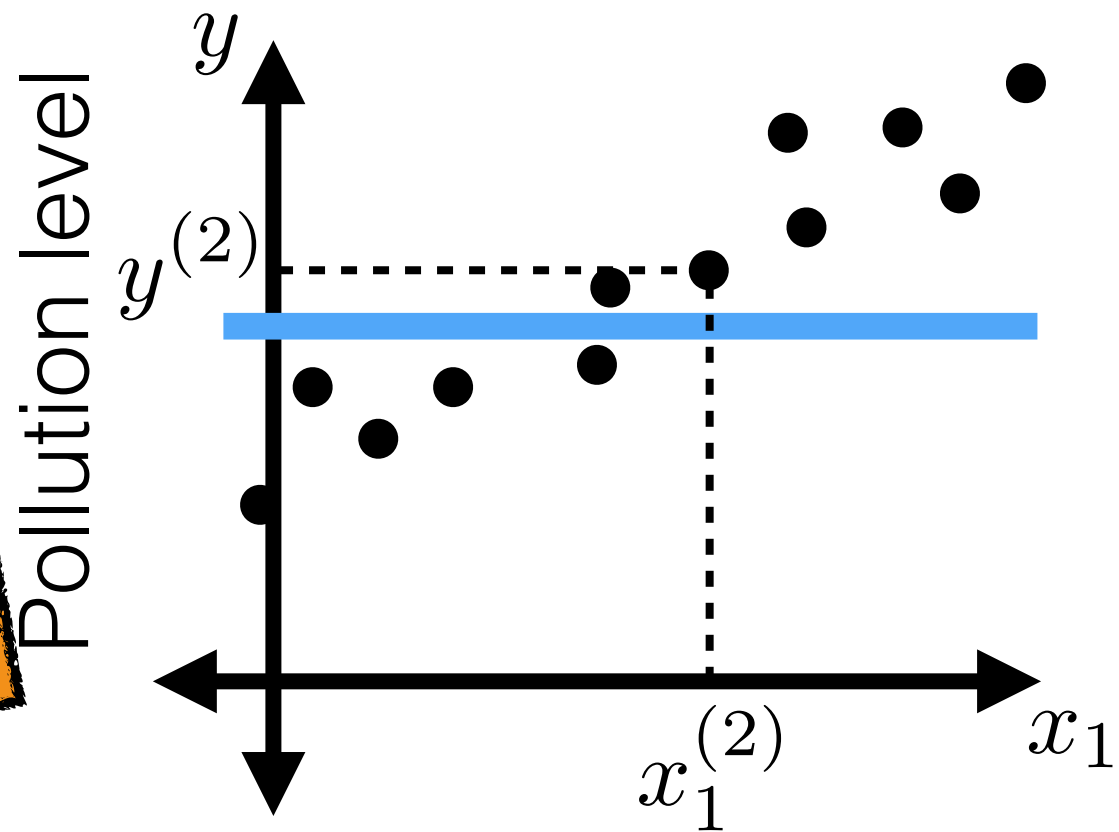
1x2, 2x1

OR

$$\begin{aligned} h(x; \theta) &= \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1) \\ &= \theta^\top x \end{aligned}$$

1x3, 3x1

Notational trick: not the same  $\theta$  &  $x$ !



# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters

- A linear reg. hypothesis when  $d \geq 1$ :

$$\begin{aligned} h(x; \theta, \theta_0) &= \theta_1 x_1 + \dots + \theta_d x_d + \theta_0 \\ &= \theta^\top x + \theta_0 \end{aligned}$$

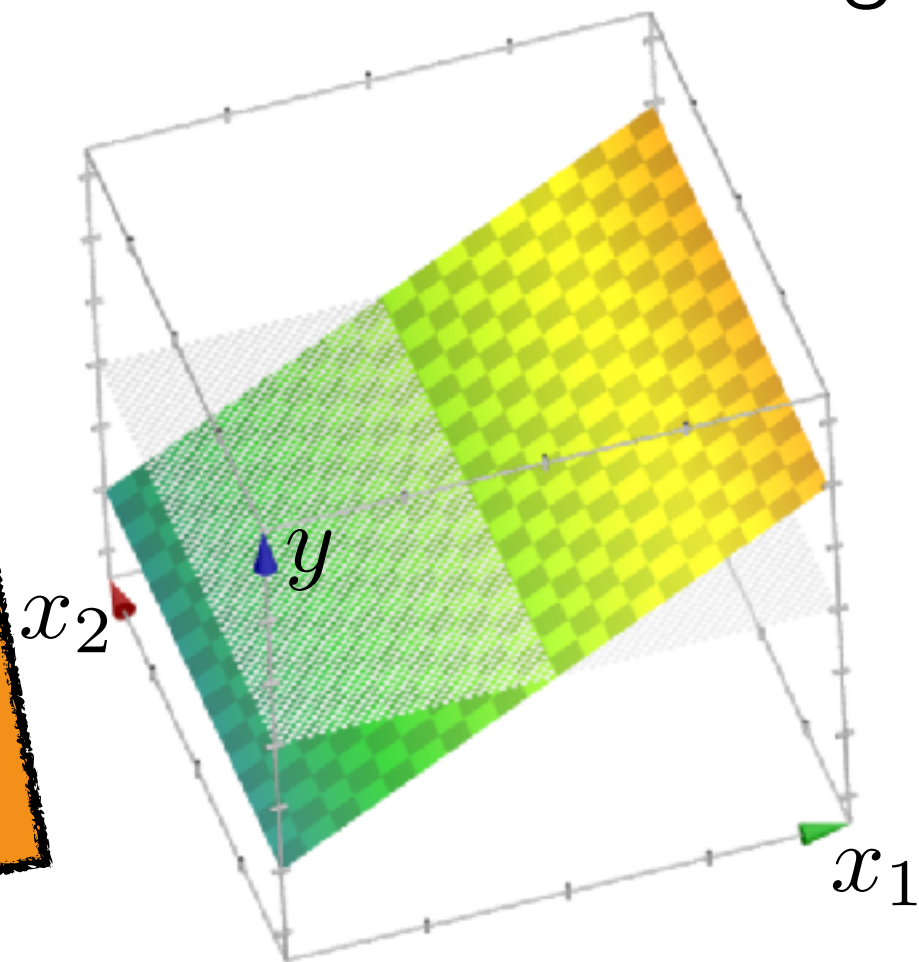
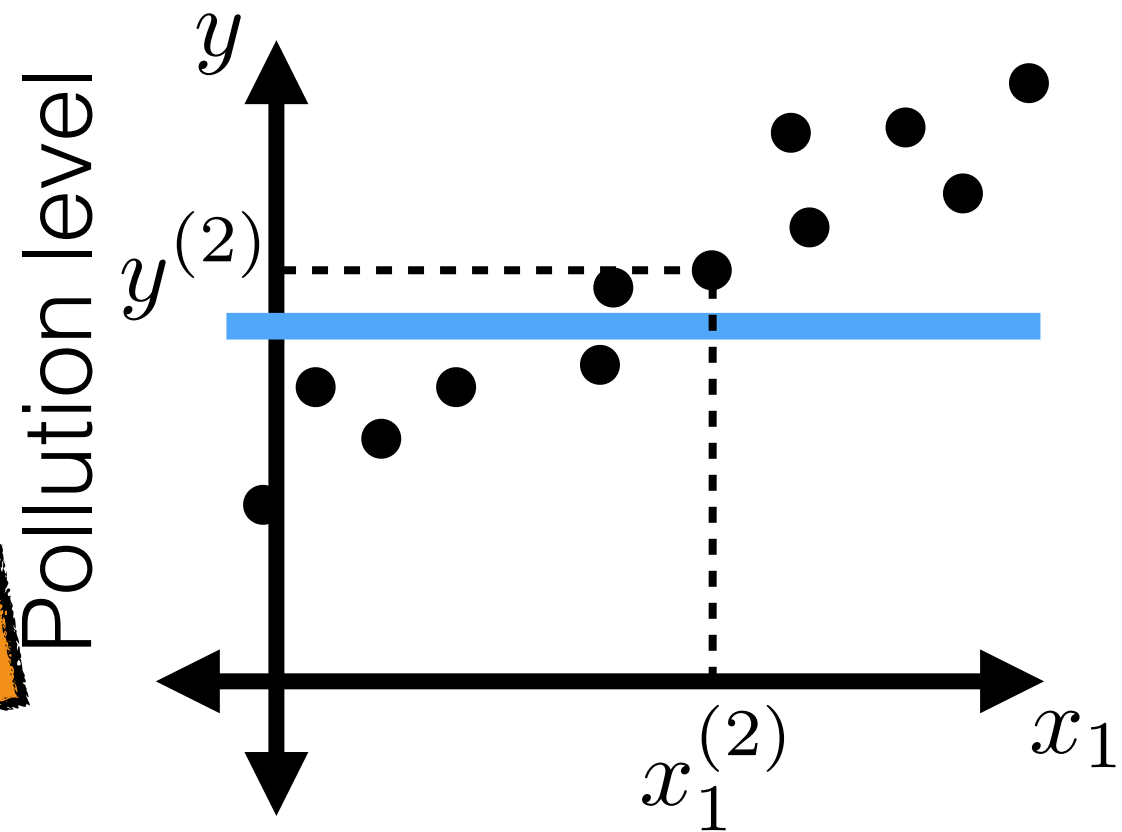
1x2, 2x1

OR

$$\begin{aligned} h(x; \theta) &= \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1) \\ &= \theta^\top x \end{aligned}$$

1x3, 3x1

Notational trick: not the same  $\theta$  &  $x$ !

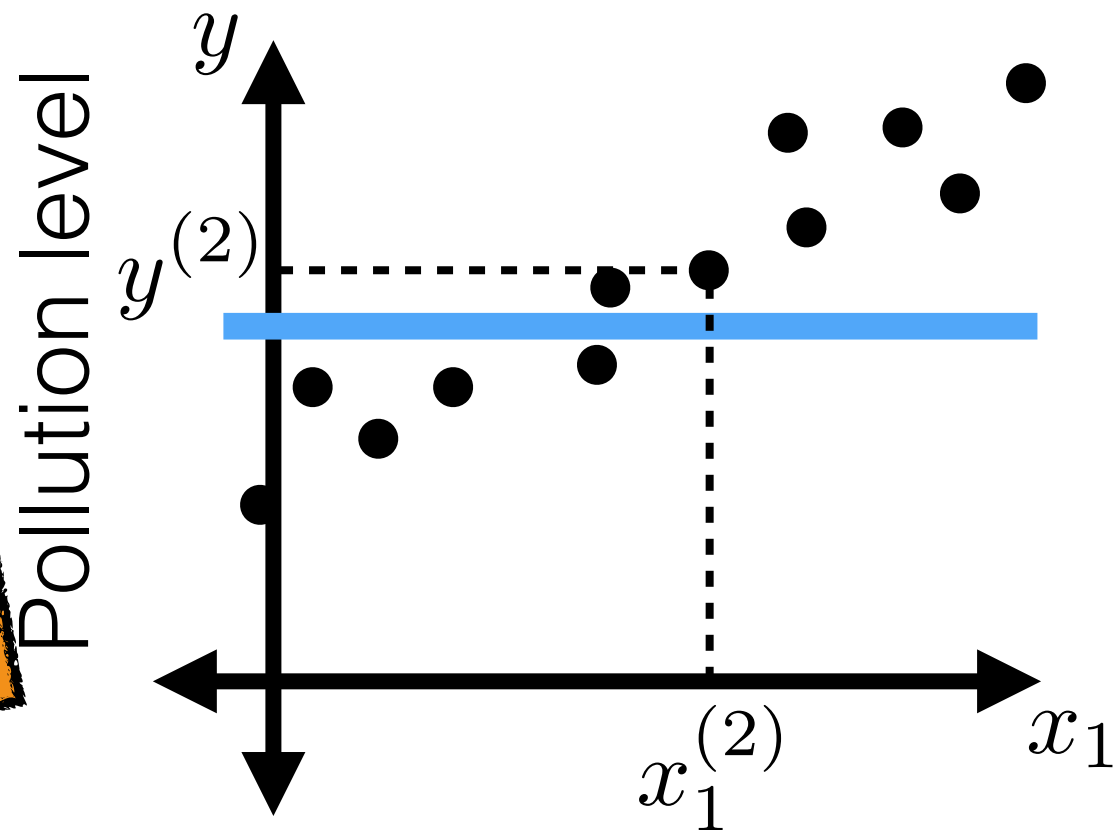


# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters



- A linear reg. hypothesis when  $d \geq 1$ :

$$h(x; \theta, \theta_0) = \theta_1 x_1 + \dots + \theta_d x_d + \theta_0$$

$$= \theta^\top x + \theta_0$$

1x2, 2x1

Hypothesis is a "hyperplane"

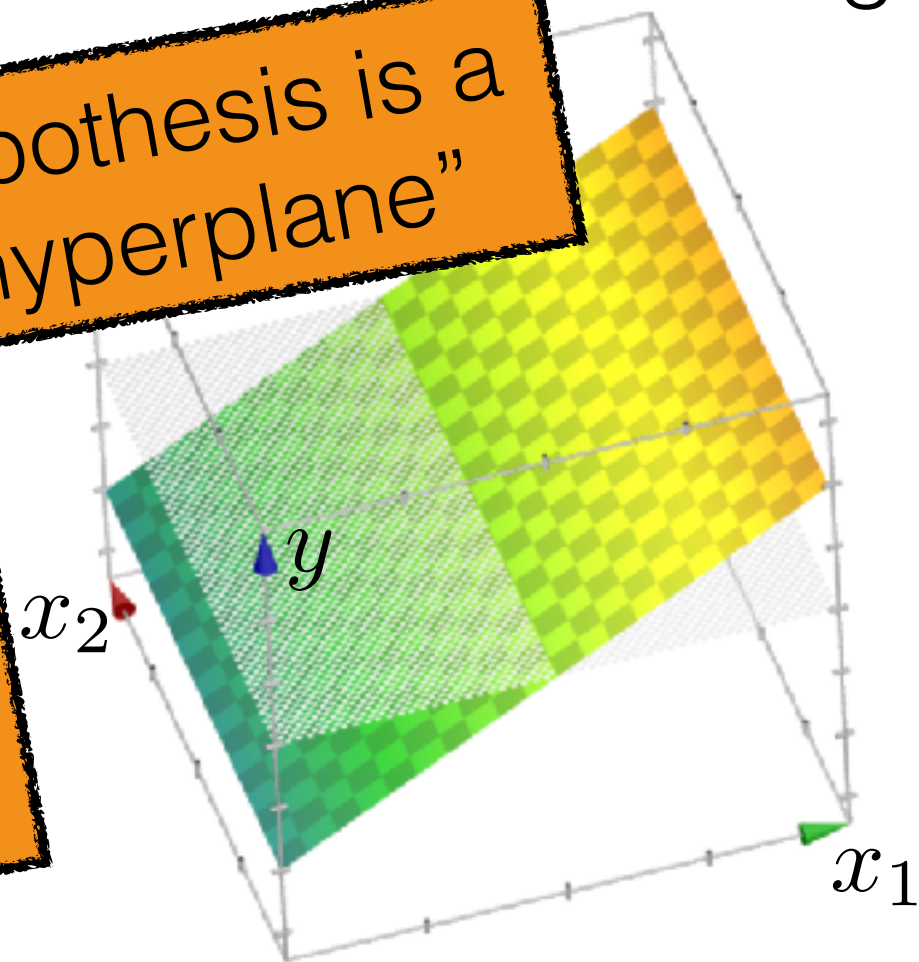
OR

$$h(x; \theta) = \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1)$$

$$= \theta^\top x$$

1x3, 3x1

Notational trick: not the same  $\theta$  &  $x$ !



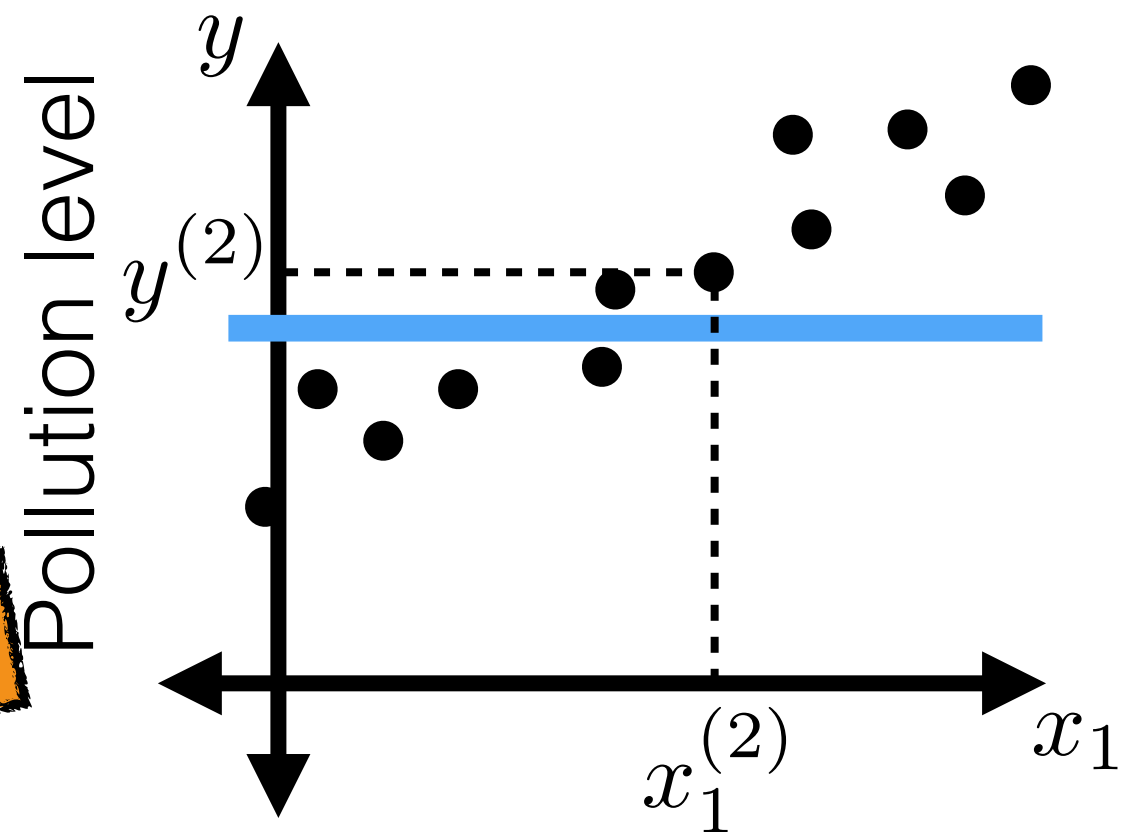


# Linear regressors

- Hypothesis class  $\mathcal{H}$ : set of  $h$ 
  - Example: all constant functions
- A linear regression hypothesis when  $d=1$ :

$$h(x; \theta, \theta_0) = \theta x + \theta_0$$

parameters



- A linear reg. hypothesis when  $d \geq 1$ :

$$h(x; \theta, \theta_0) = \theta_1 x_1 + \dots + \theta_d x_d + \theta_0$$

$$= \theta^\top x + \theta_0$$

1x2, 2x1

Hypothesis is a "hyperplane"

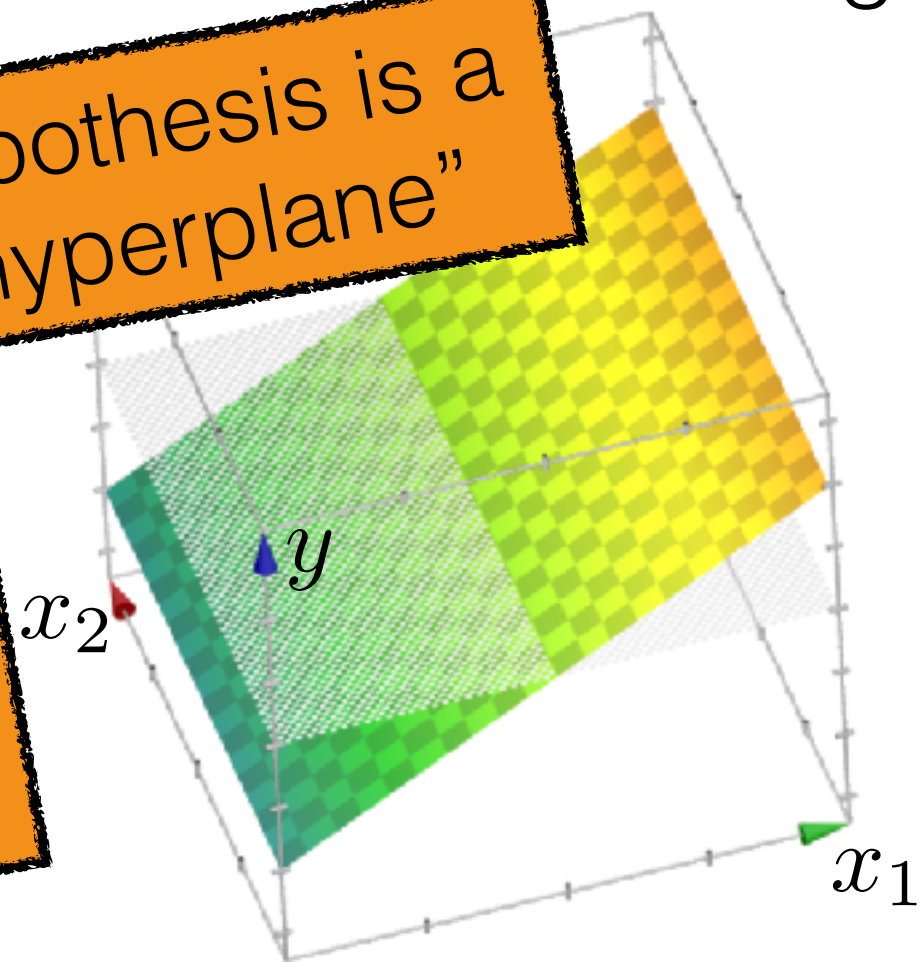
OR

$$h(x; \theta) = \theta_1 x_1 + \dots + \theta_d x_d + (\theta_0)(1)$$

$$= \theta^\top x$$

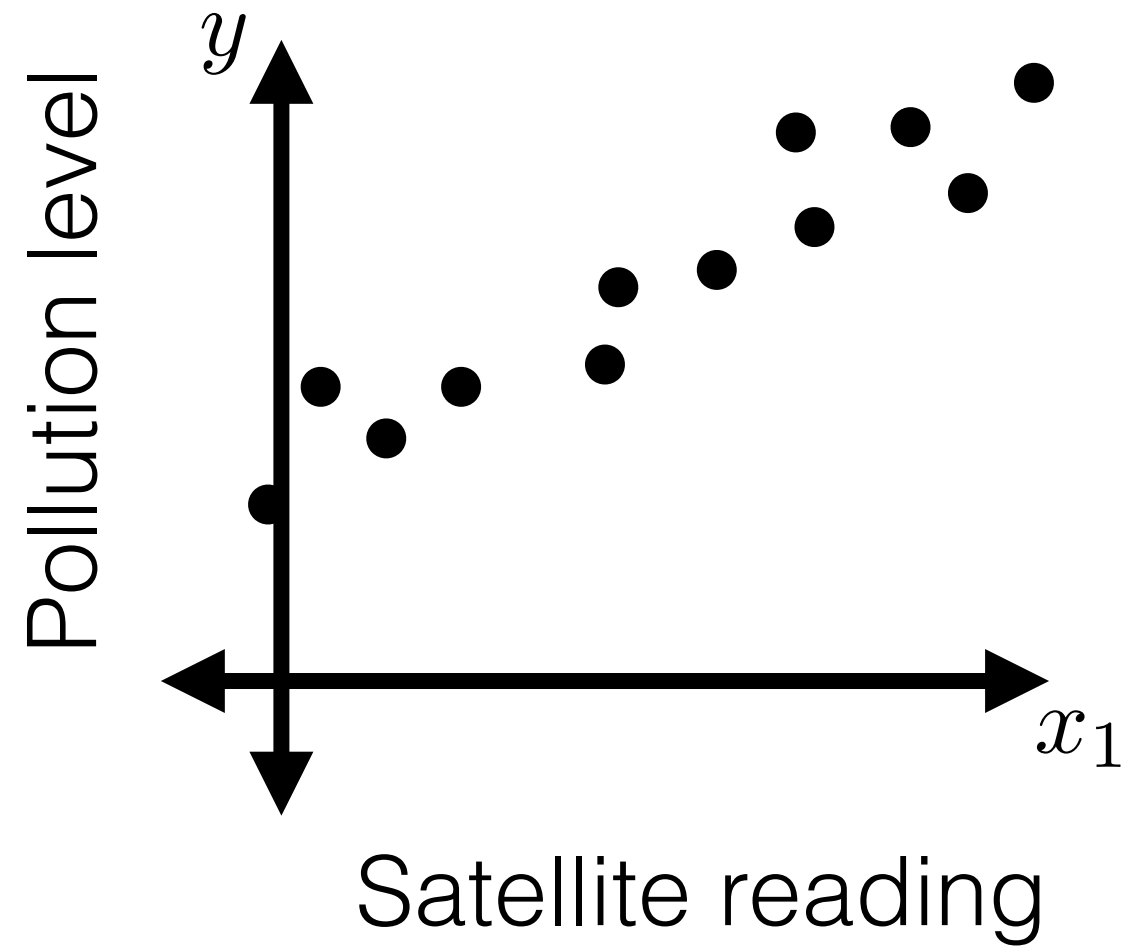
1x3, 3x1

Notational trick: not the same  $\theta$  &  $x$ !



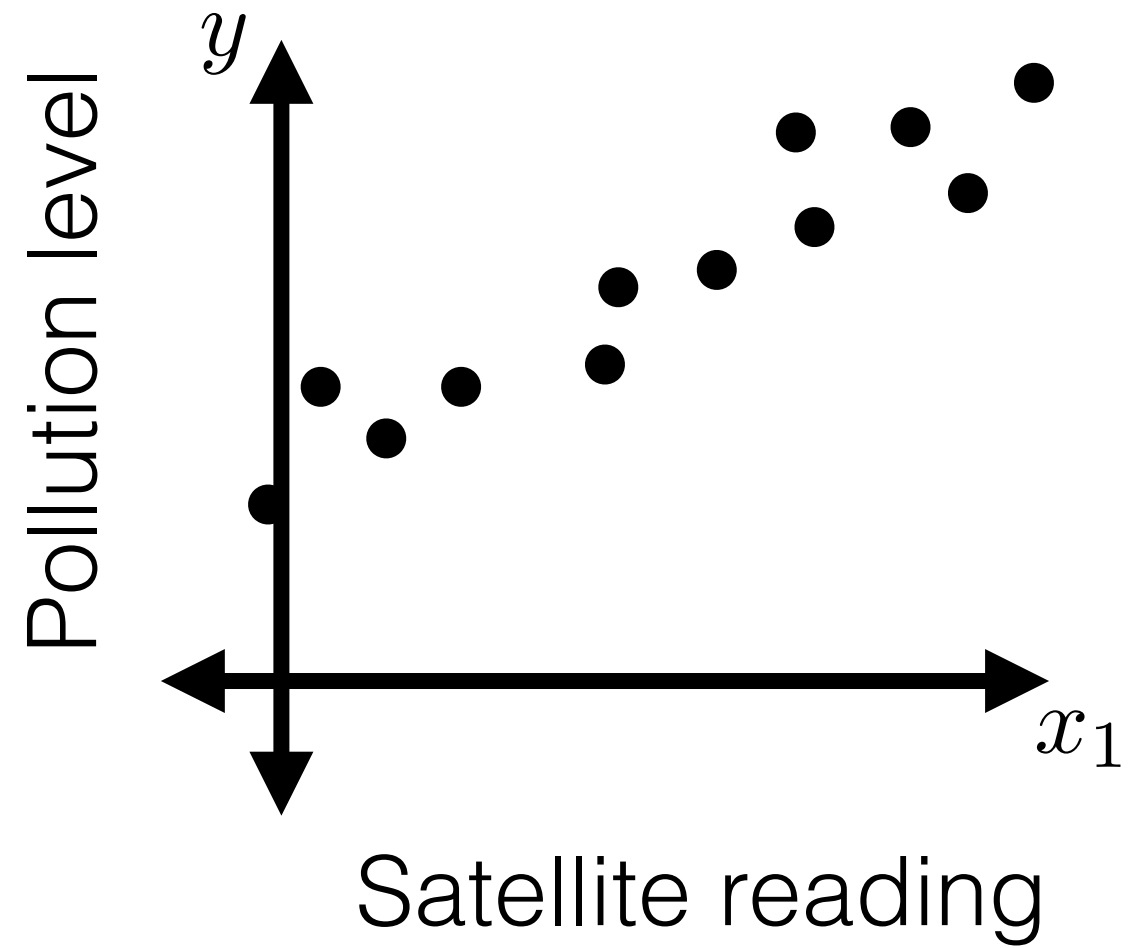
- Our hypothesis class in linear regression will be the set of all such  $h$

# How good is a regression hypothesis?



# How good is a regression hypothesis?

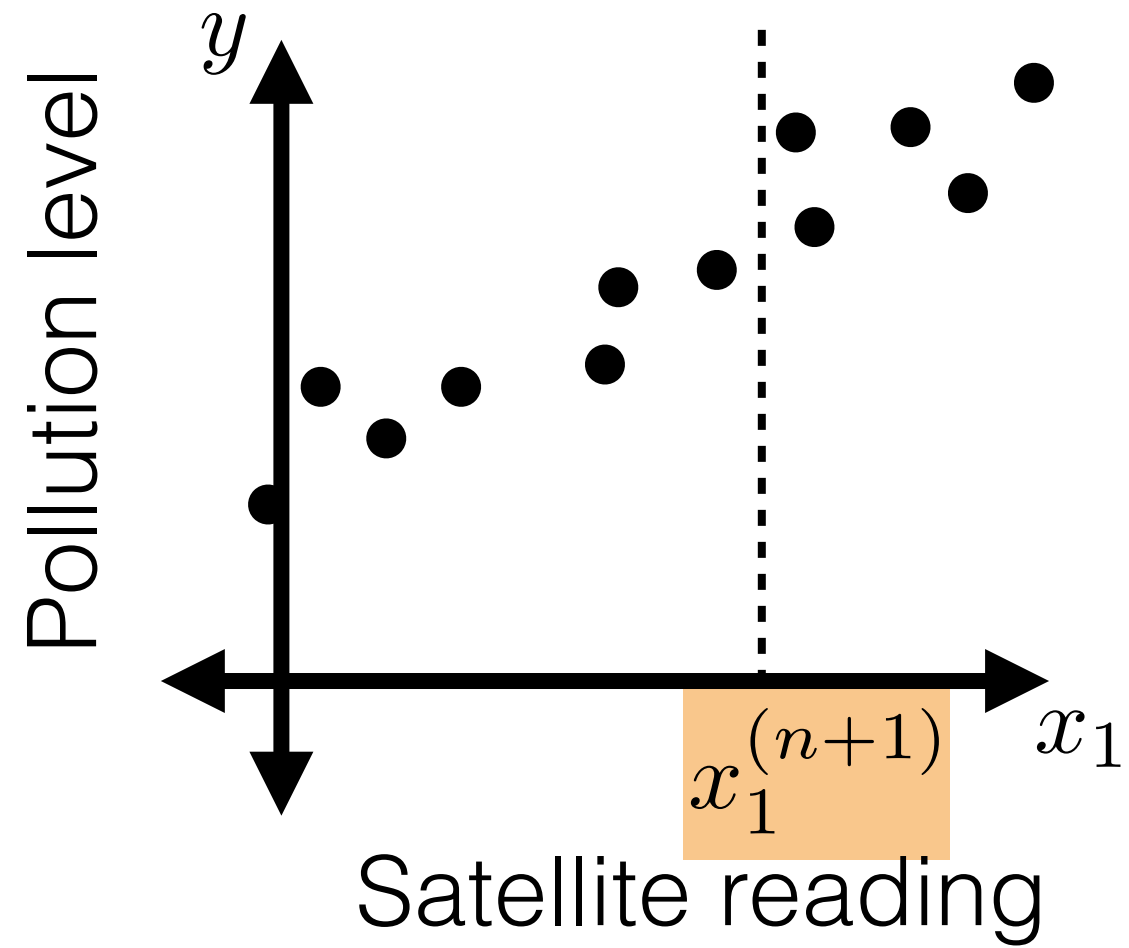
- Should predict well on future data





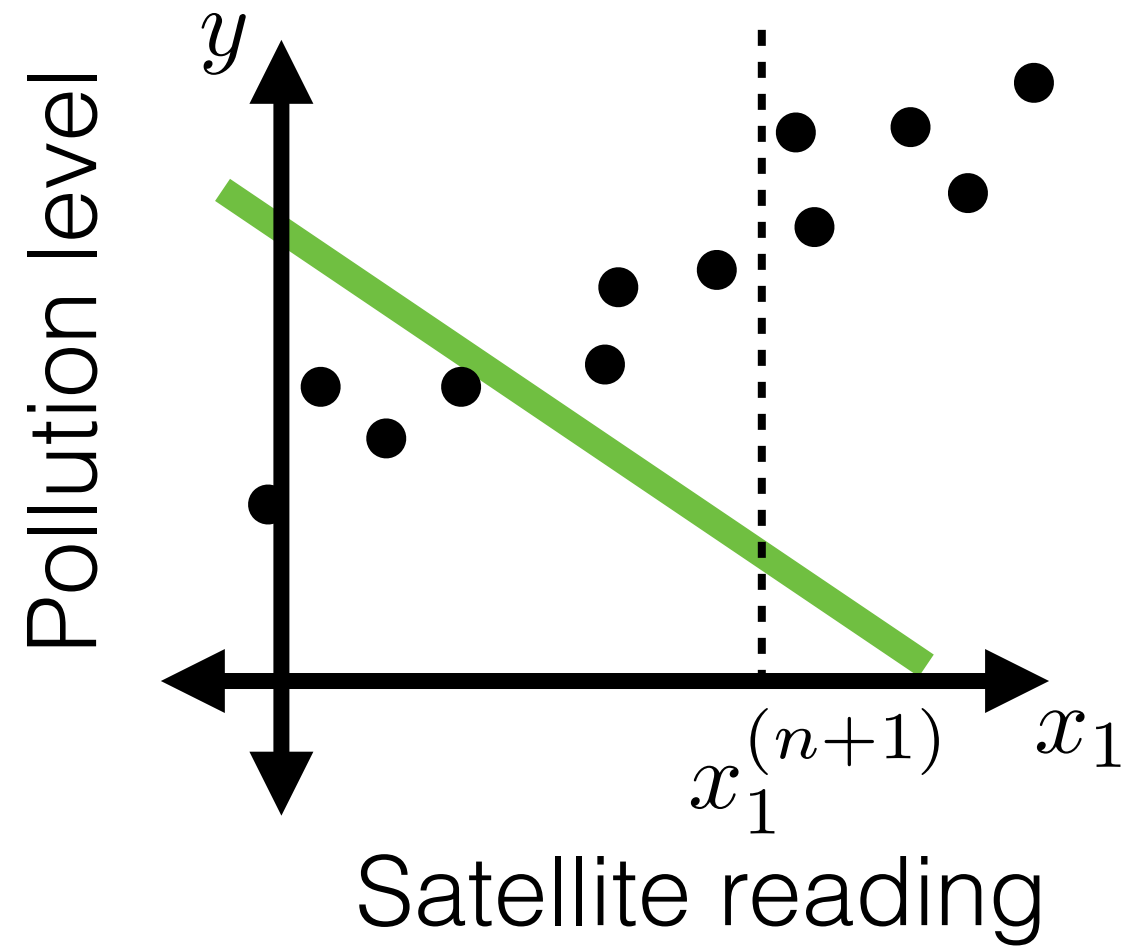
# How good is a regression hypothesis?

- Should predict well on future data



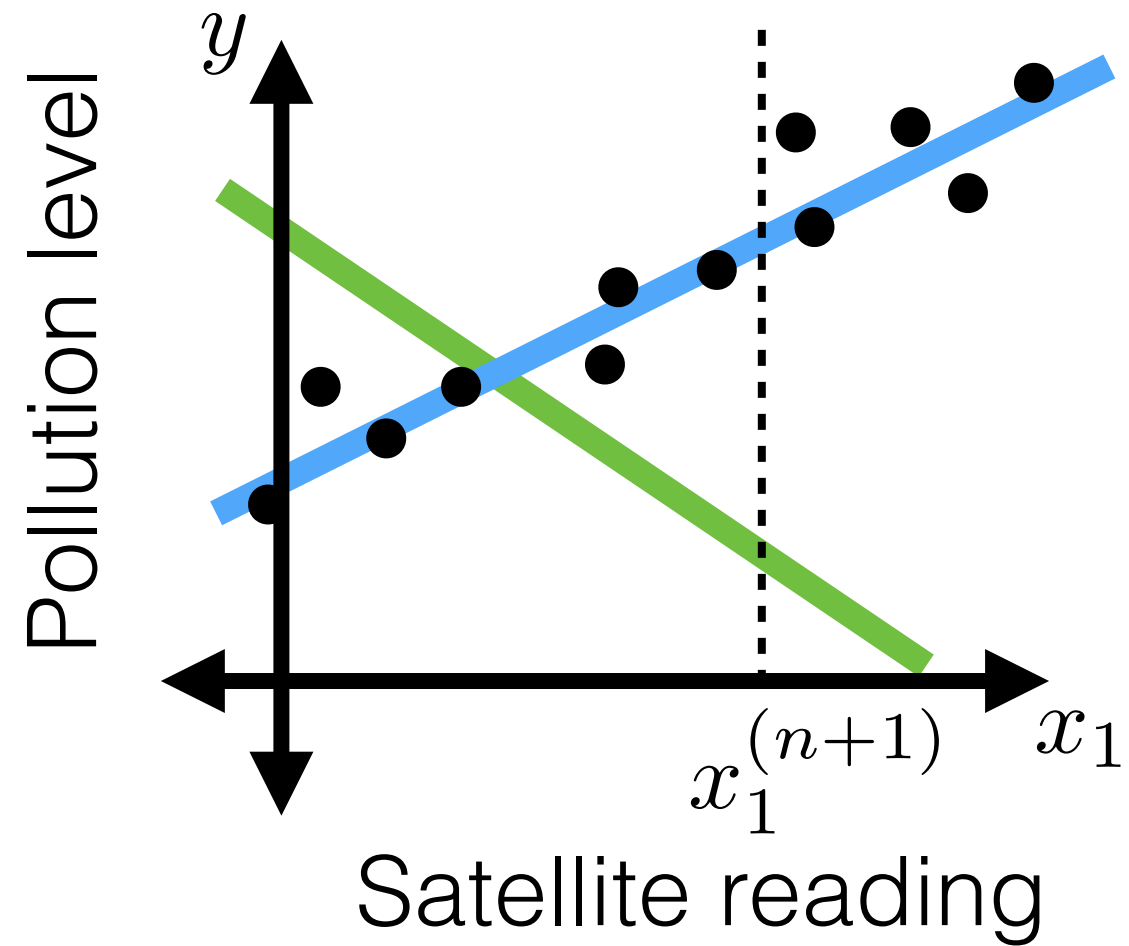
# How good is a regression hypothesis?

- Should predict well on future data



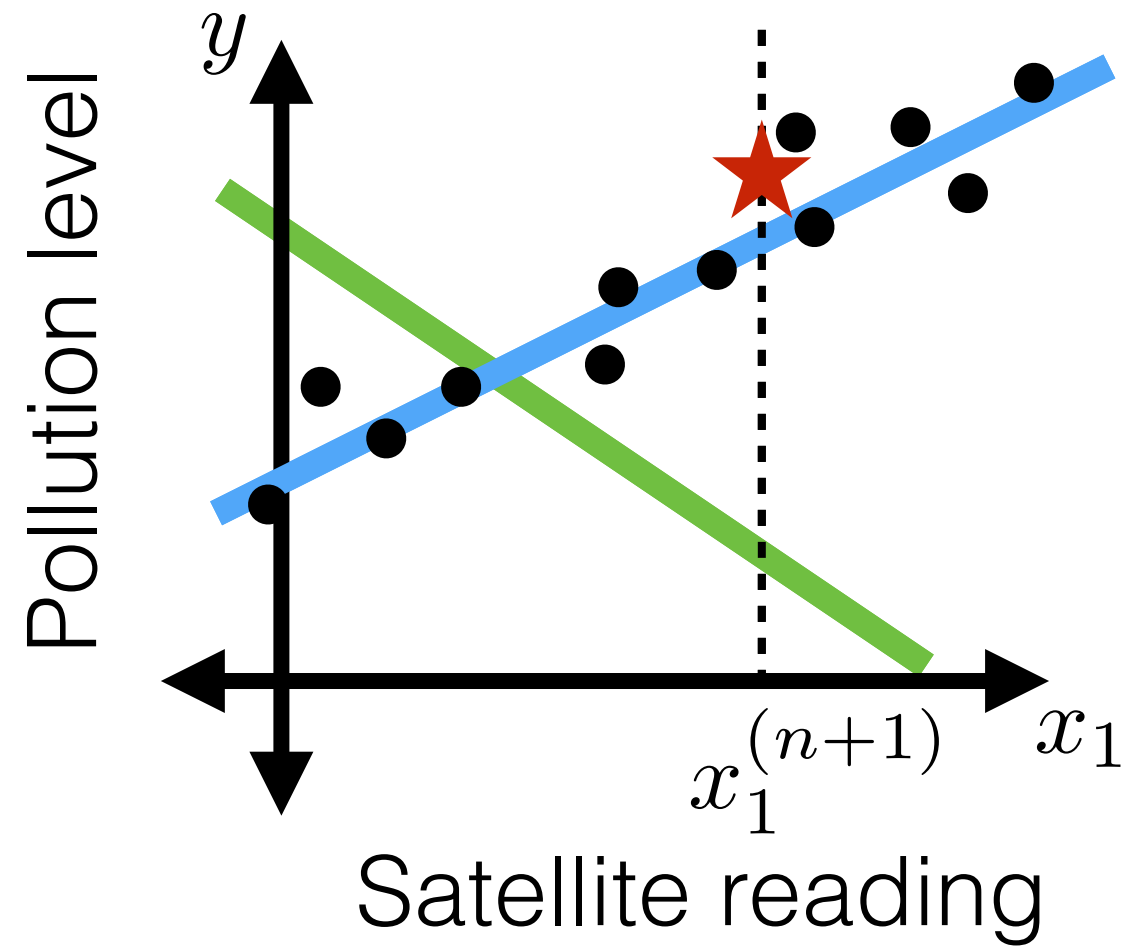
# How good is a regression hypothesis?

- Should predict well on future data



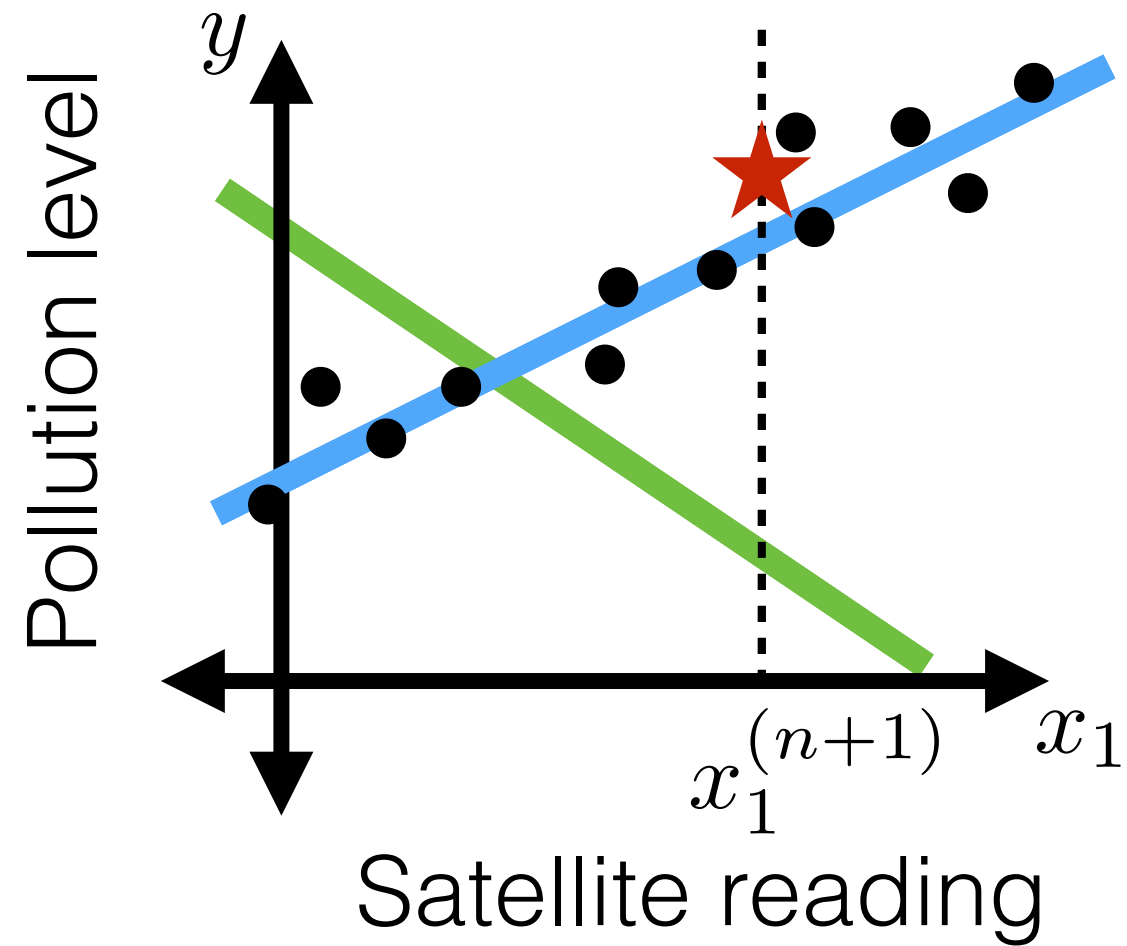
# How good is a regression hypothesis?

- Should predict well on future data



# How good is a regression hypothesis?

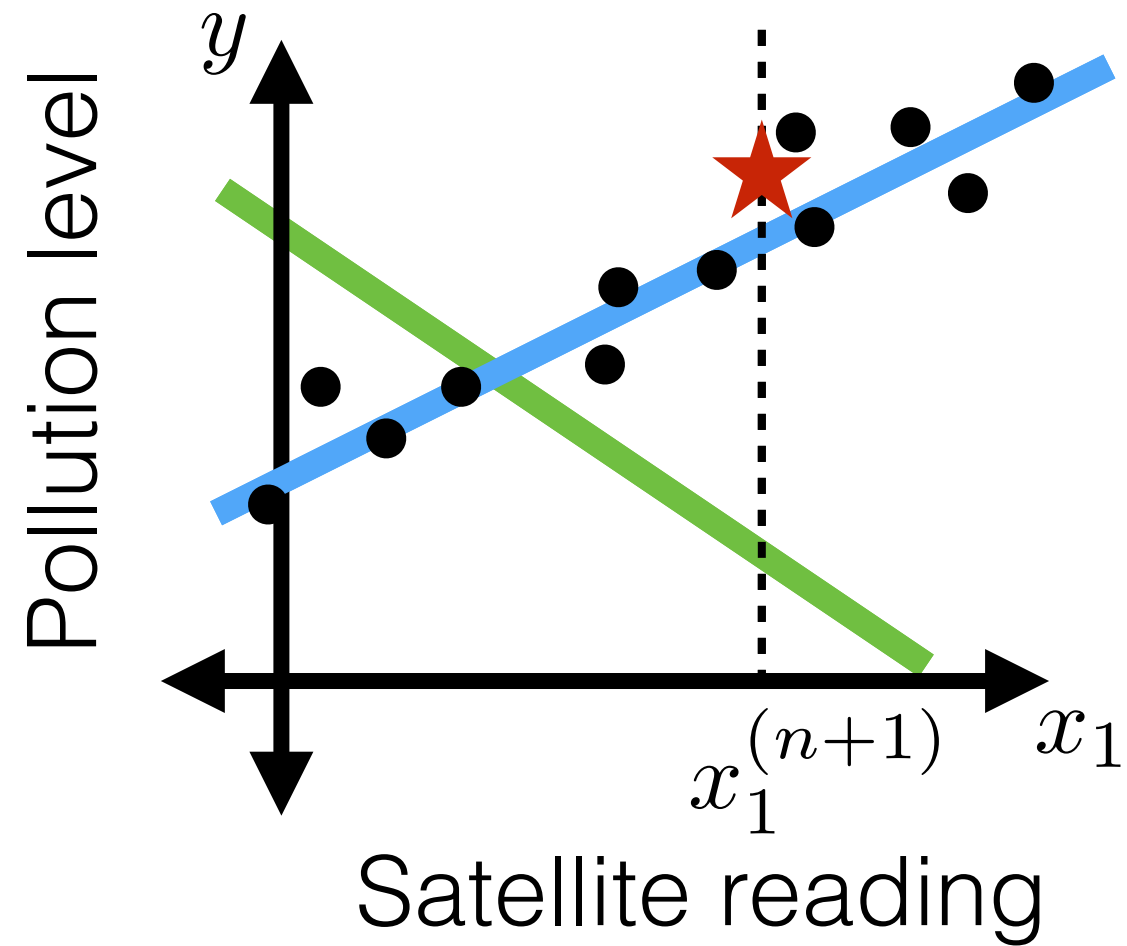
- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual



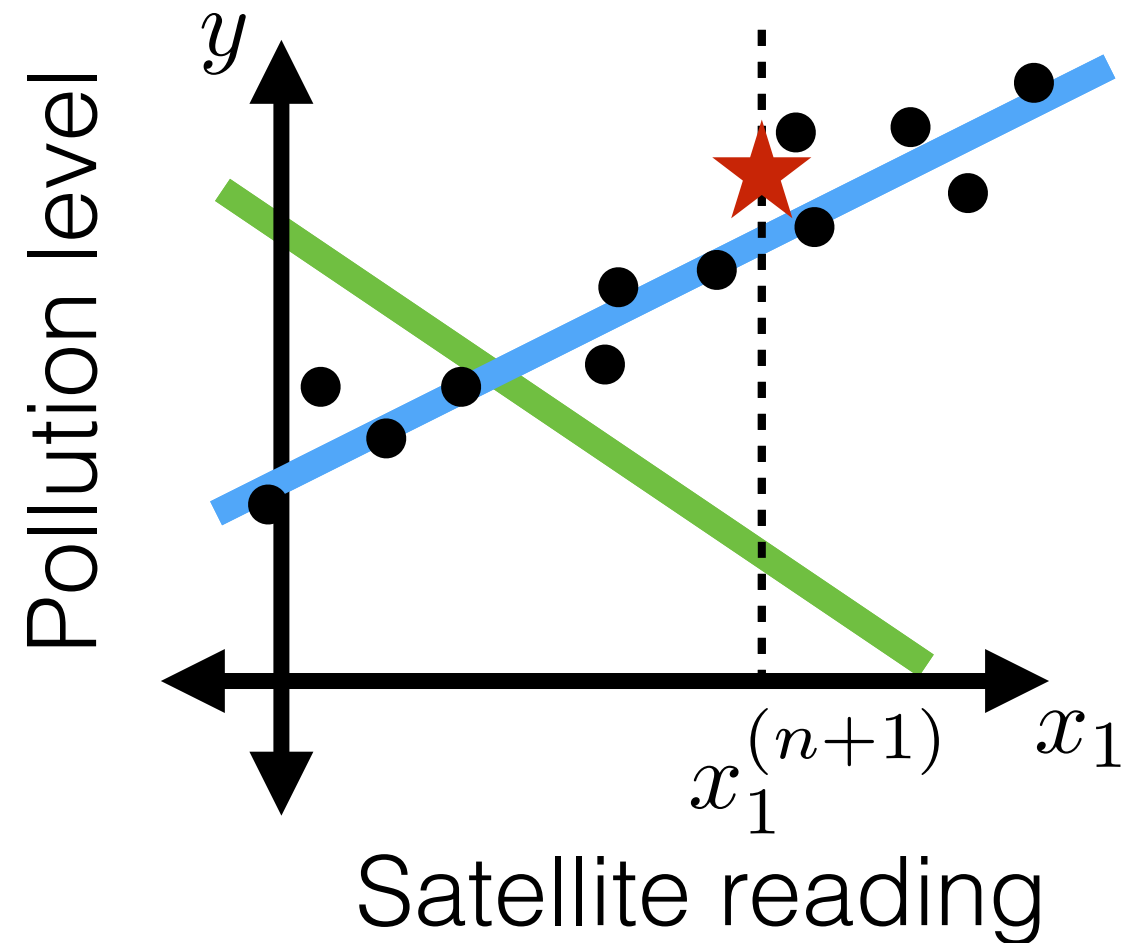
# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

- Ex: squared loss

$$L(g, a) = (g - a)^2$$





# How good is a regression hypothesis?

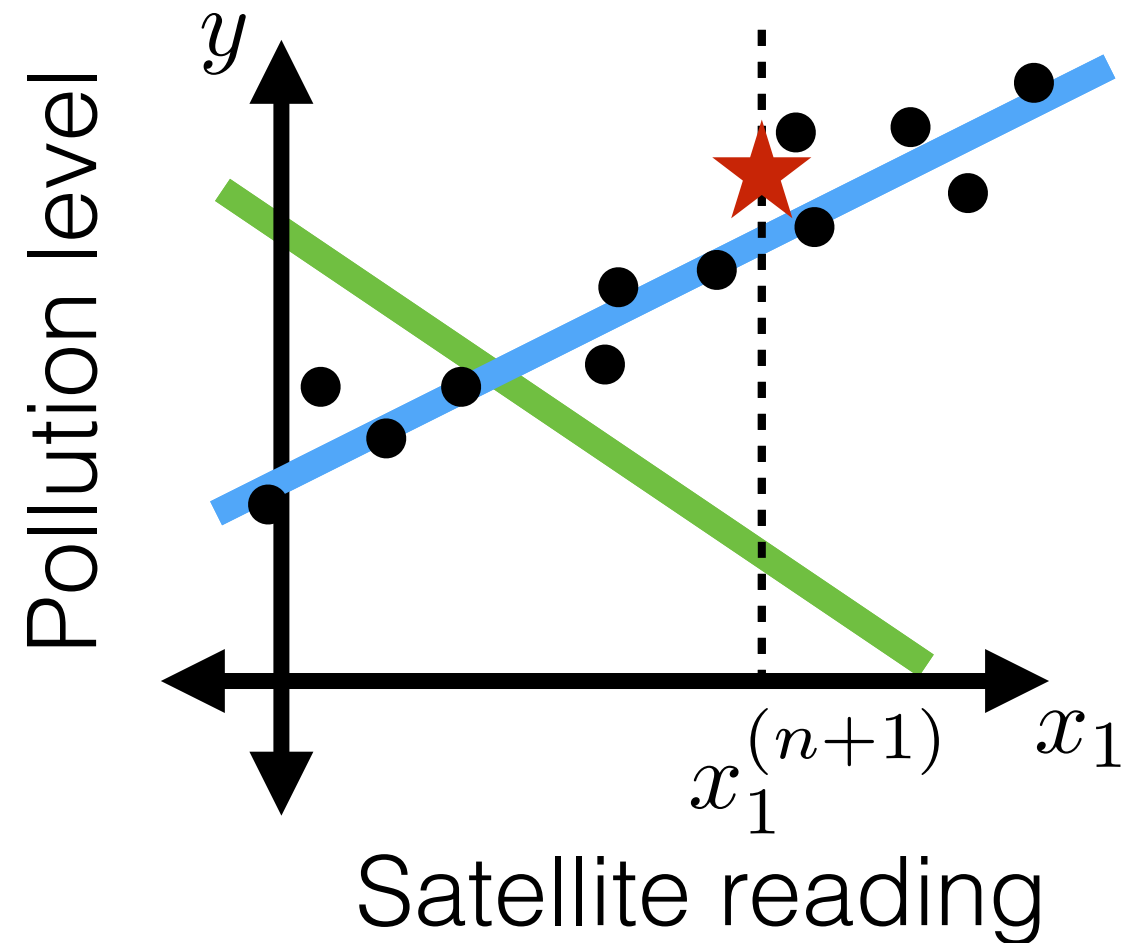
- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

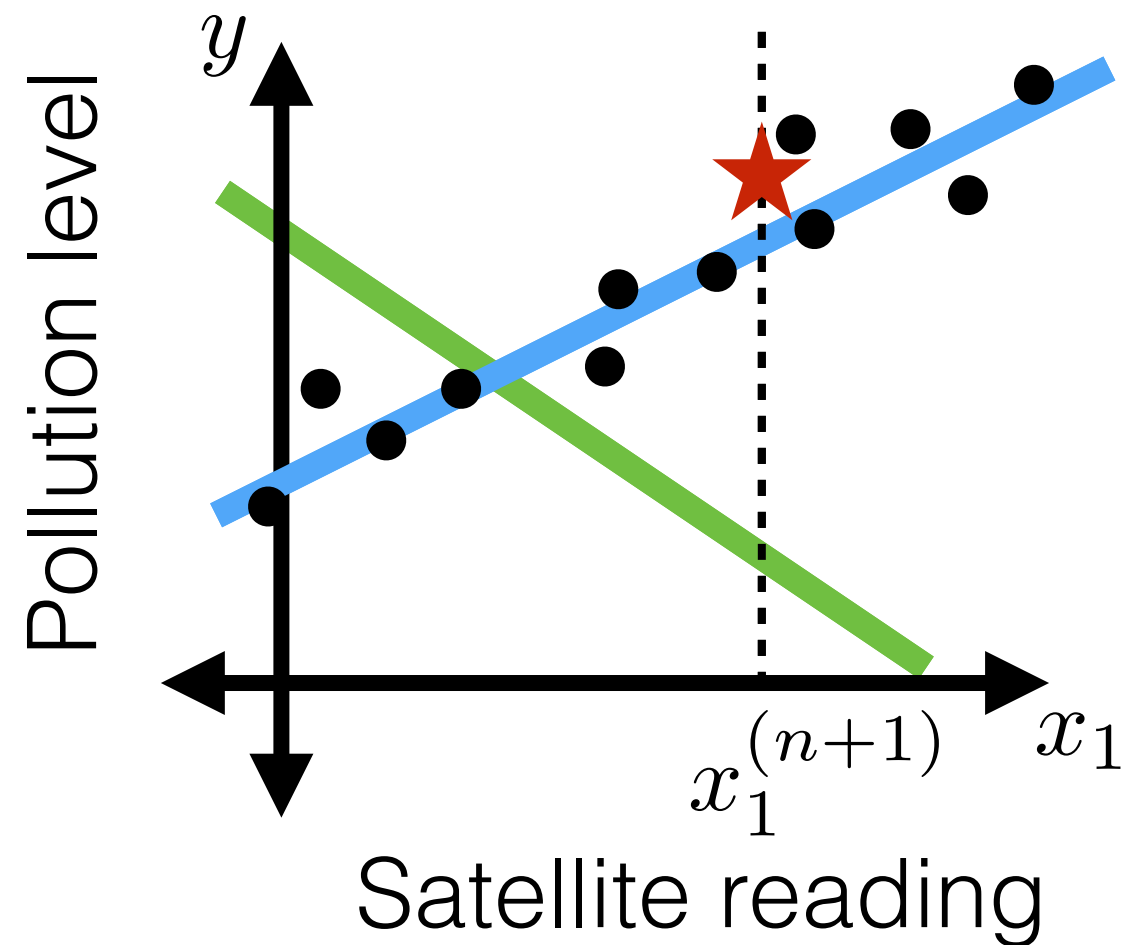
$g$ : guess,  
 $a$ : actual

- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

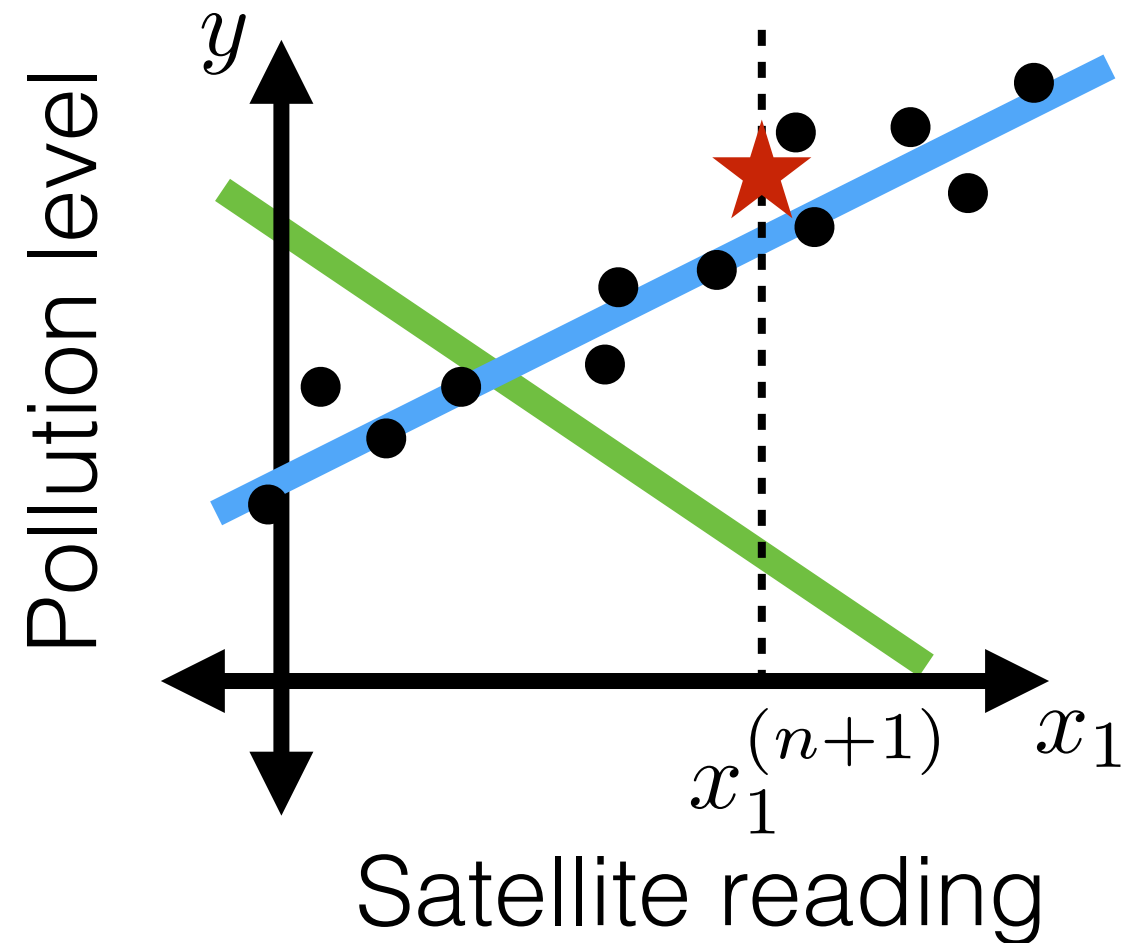
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

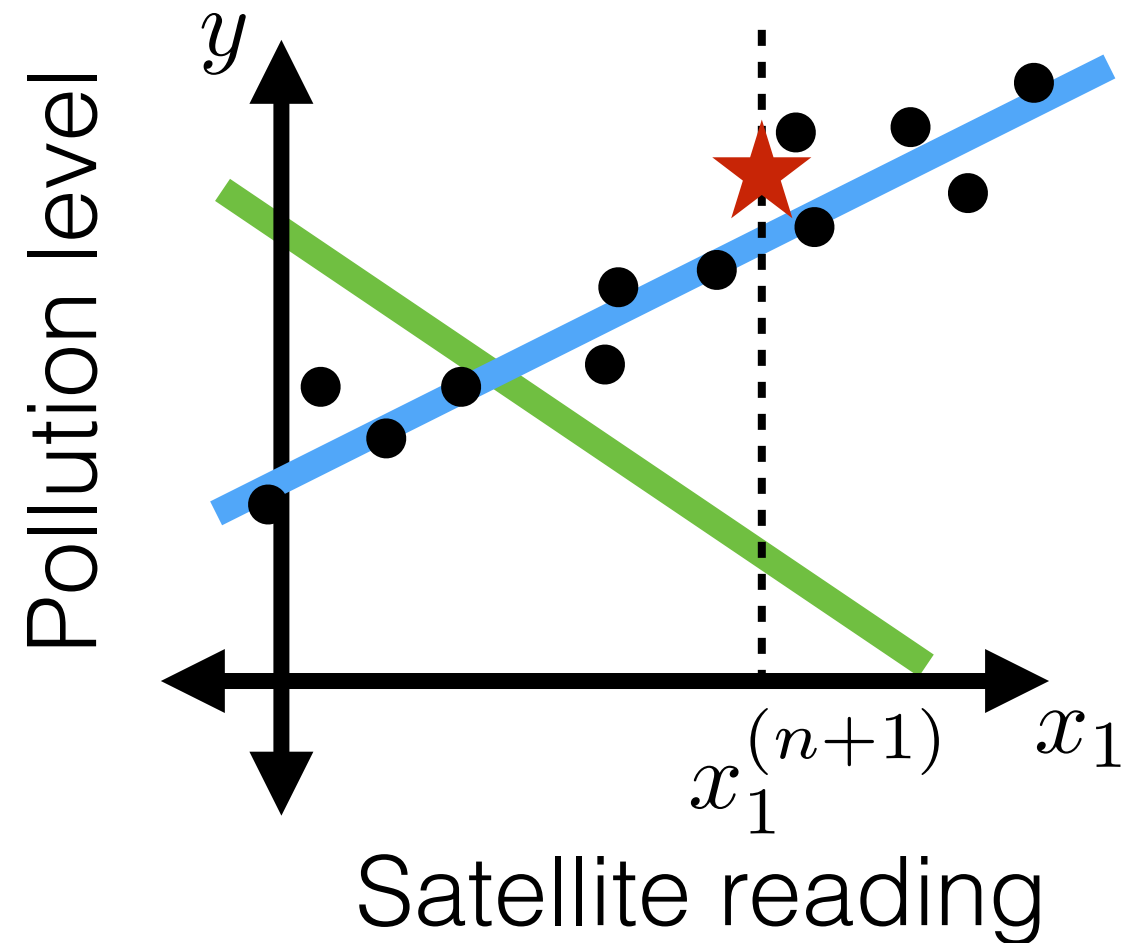
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

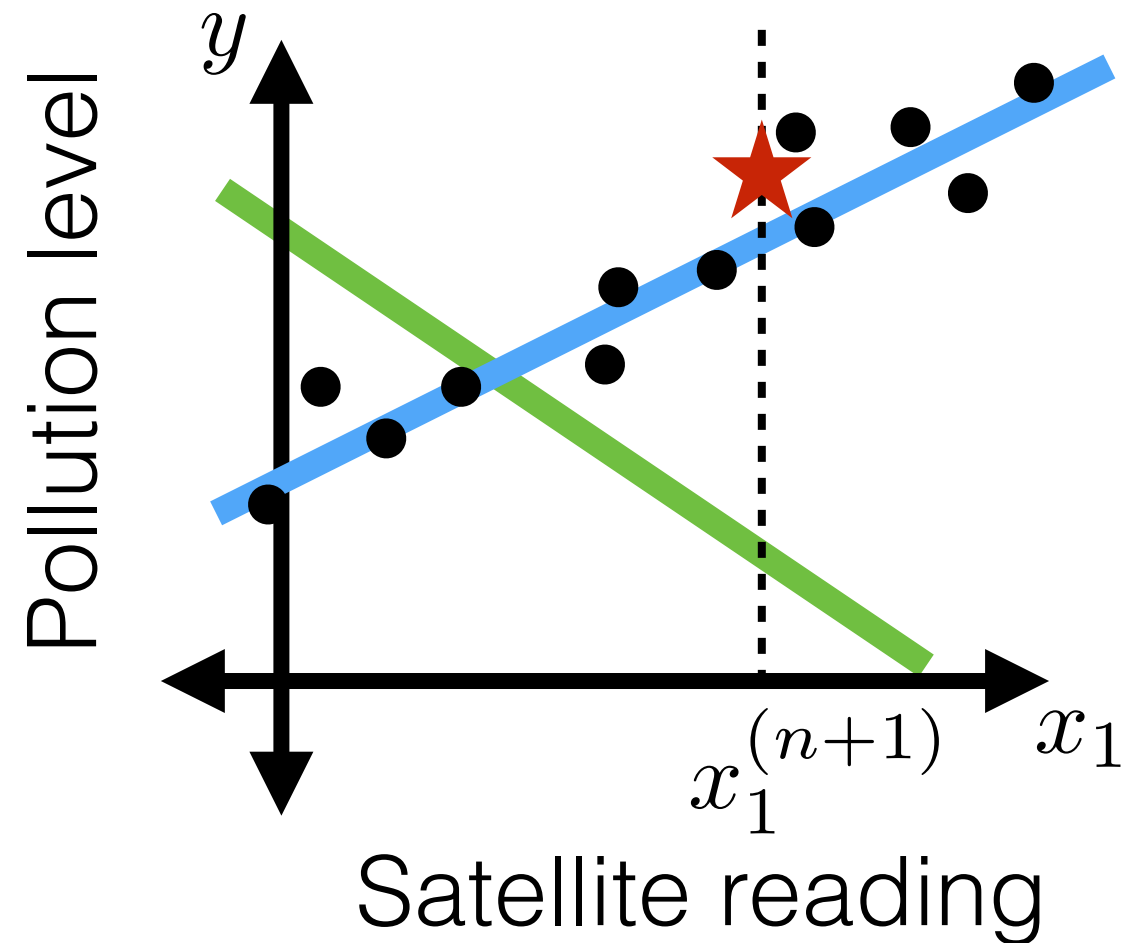
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

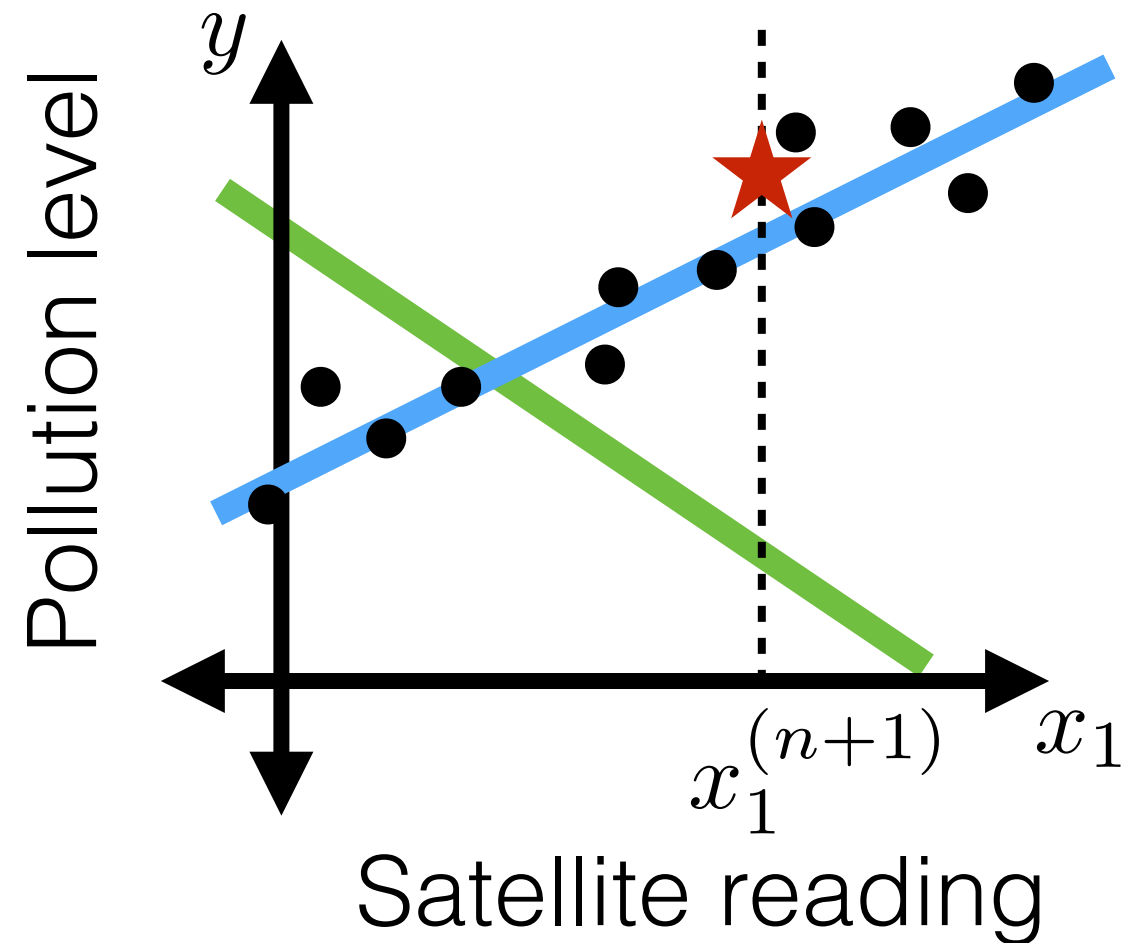
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$





# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

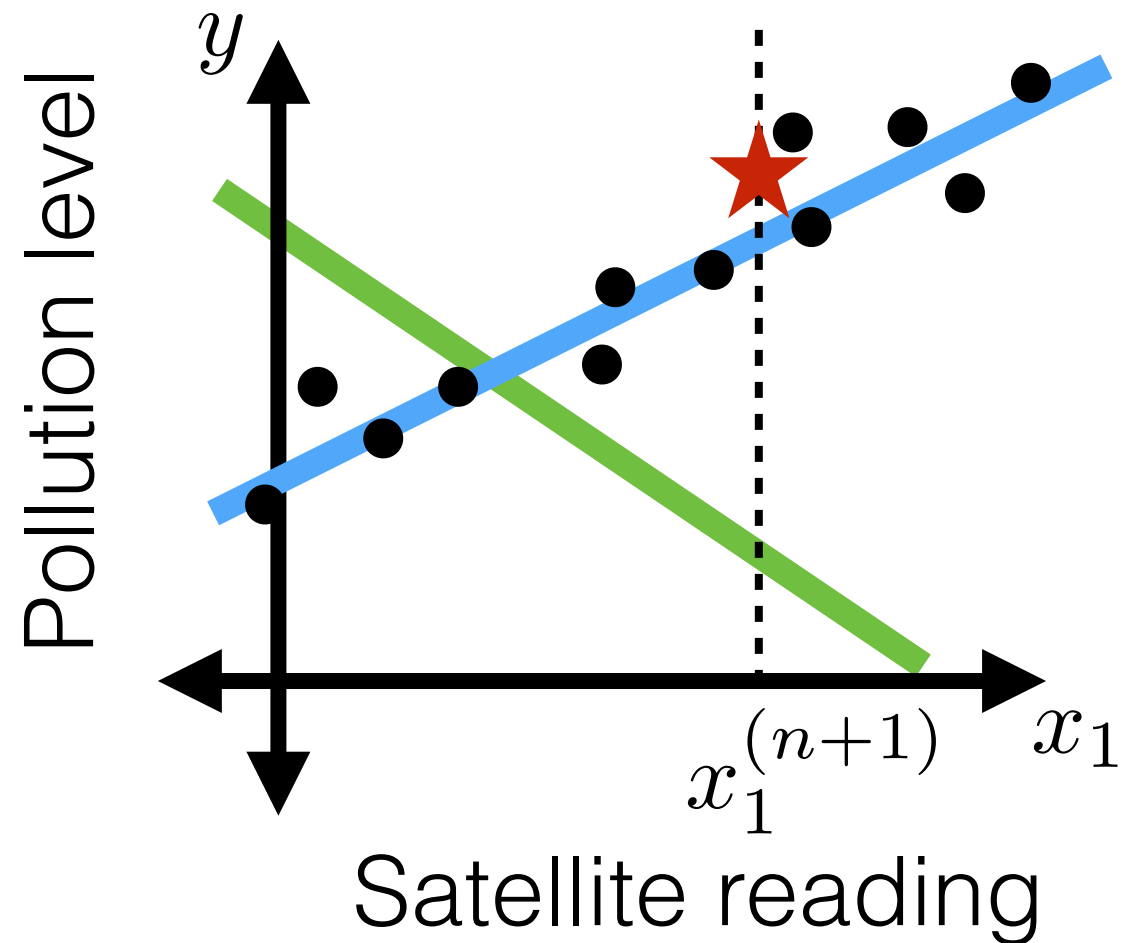
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

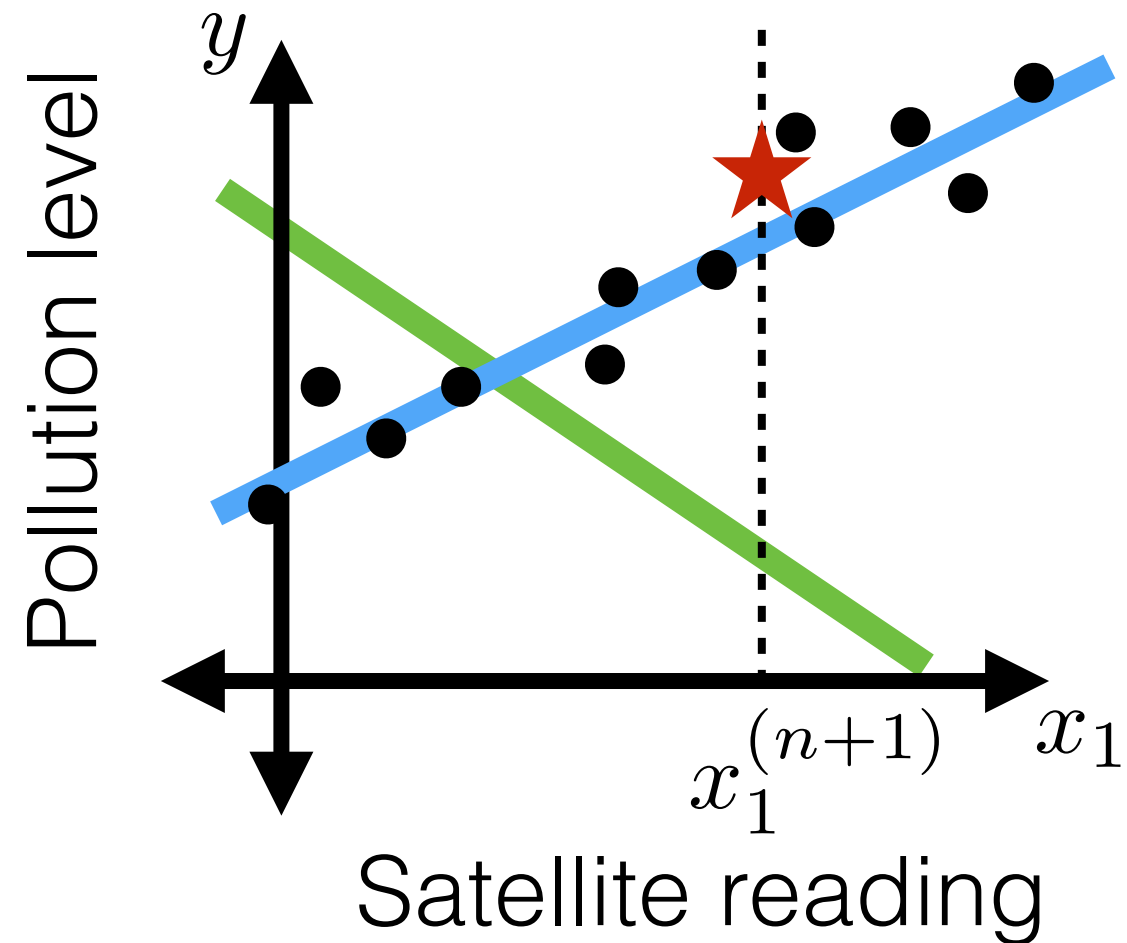
- Ex: squared loss

$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

- Ex: squared loss

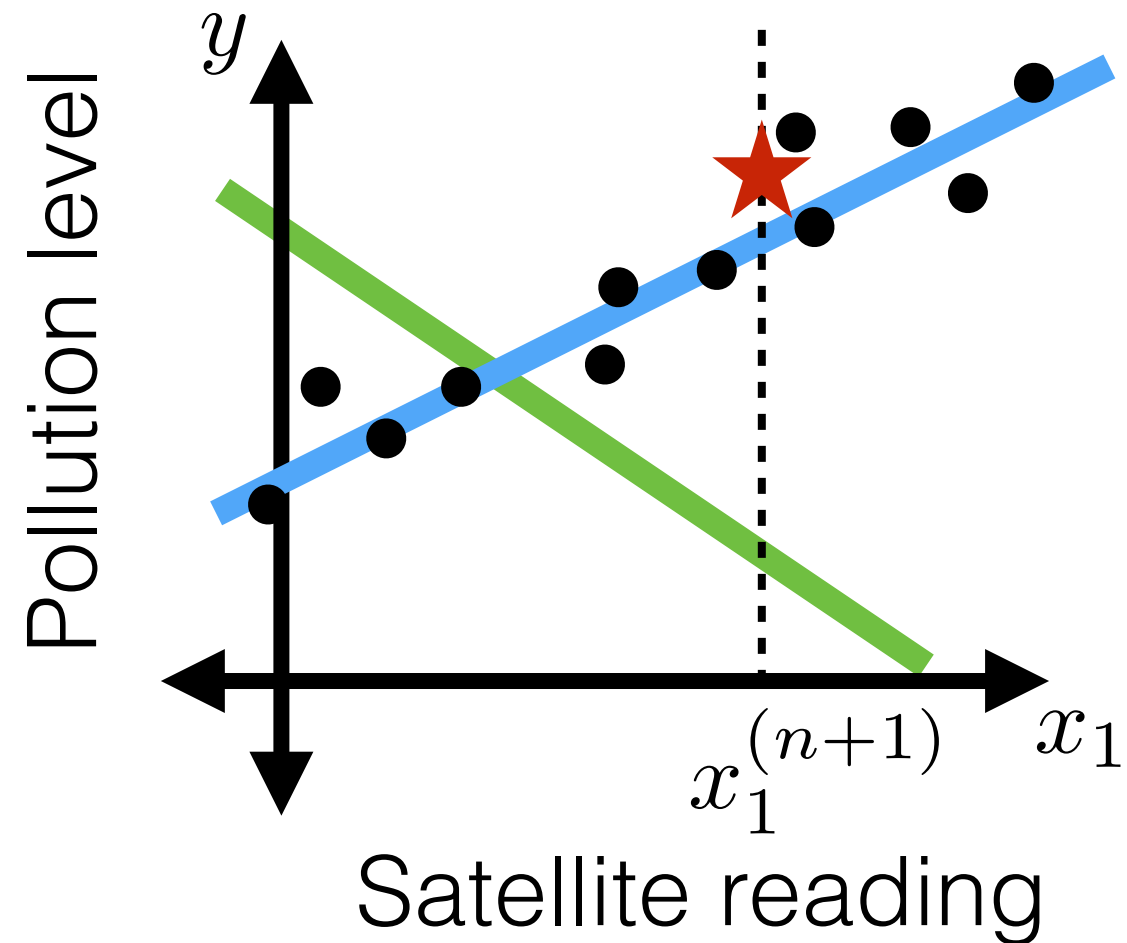
$$L(g, a) = (g - a)^2$$

- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$

- Training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$



# How good is a regression hypothesis?

- Should predict well on future data
- How good is a regressor at one point? Loss  $L(g, a)$

$g$ : guess,  
 $a$ : actual

- Ex: squared loss

$$L(g, a) = (g - a)^2$$

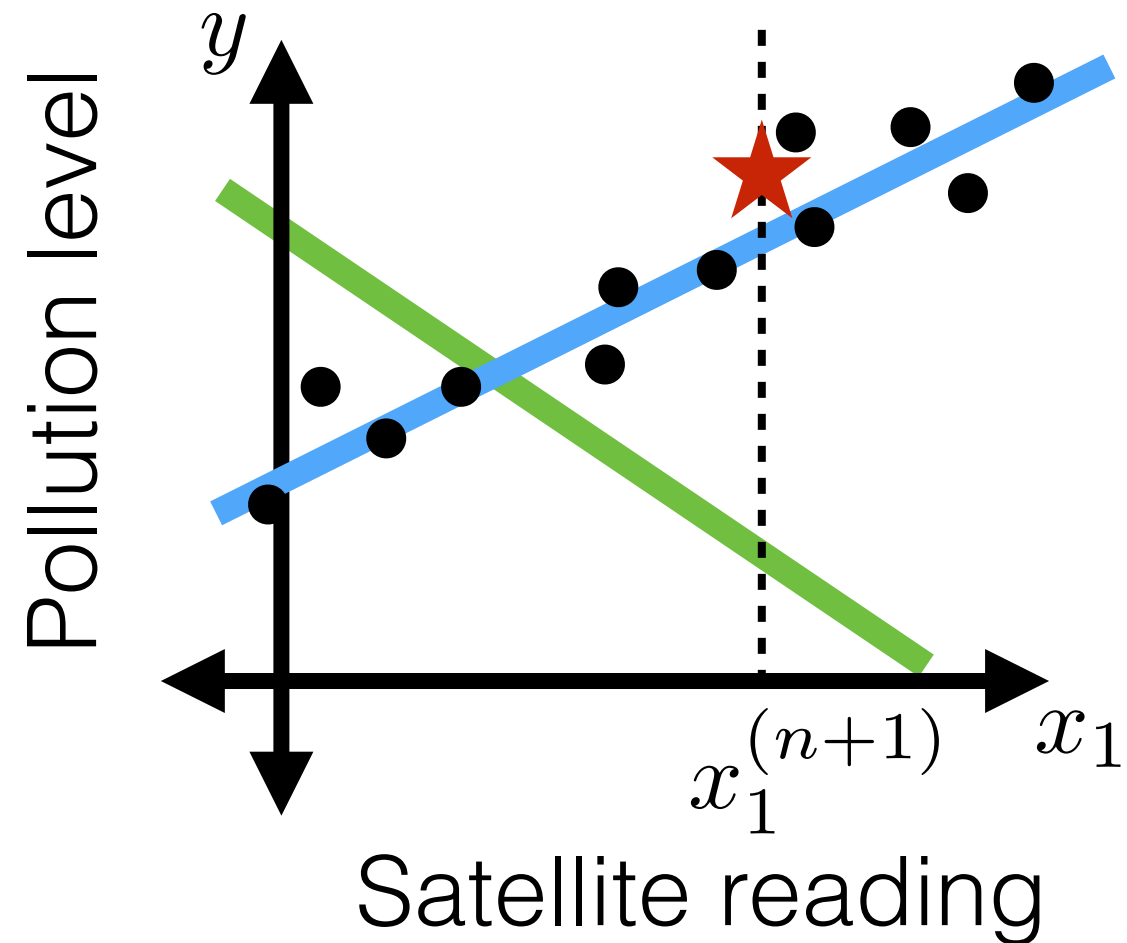
- Example: asymmetric loss

$$L(g, a) = \begin{cases} (g - a)^2 & \text{if } g > a \\ 2(g - a)^2 & \text{if } g \leq a \end{cases}$$

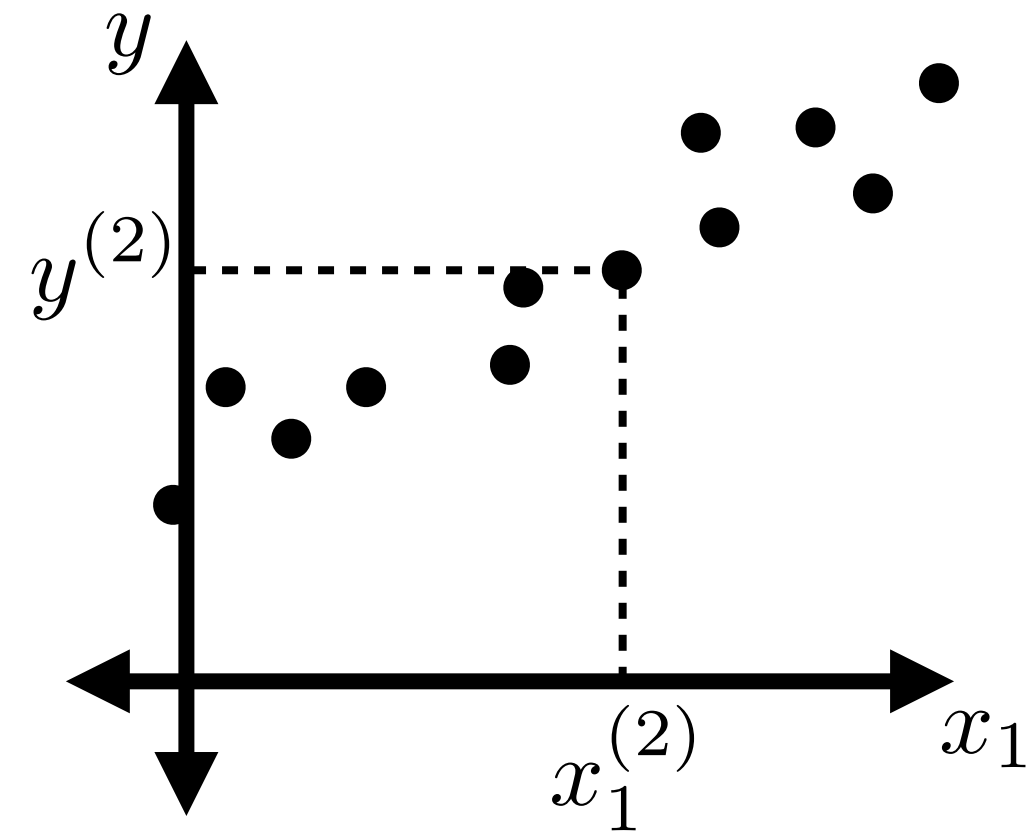
- Test error ( $n'$  new points):  $\mathcal{E}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^{(i)}), y^{(i)})$

- Training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- One idea: prefer  $h$  to  $\tilde{h}$  if  $\mathcal{E}_n(h) < \mathcal{E}_n(\tilde{h})$

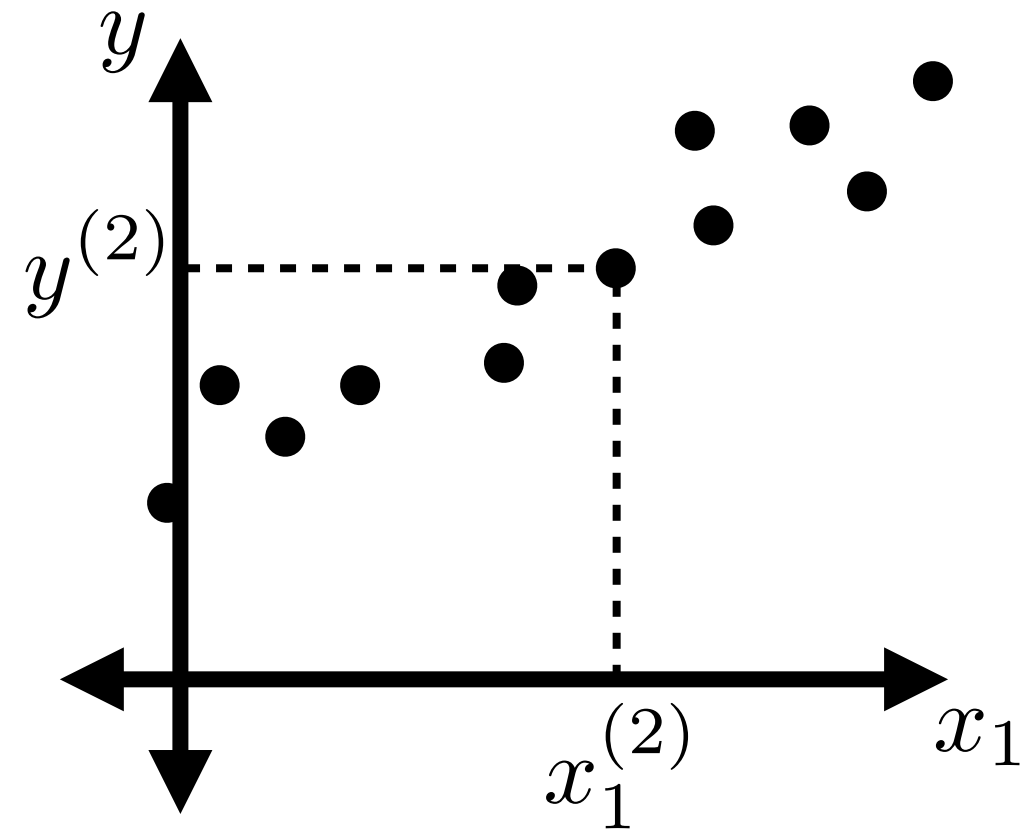


# Learning a regressor



# Learning a regressor

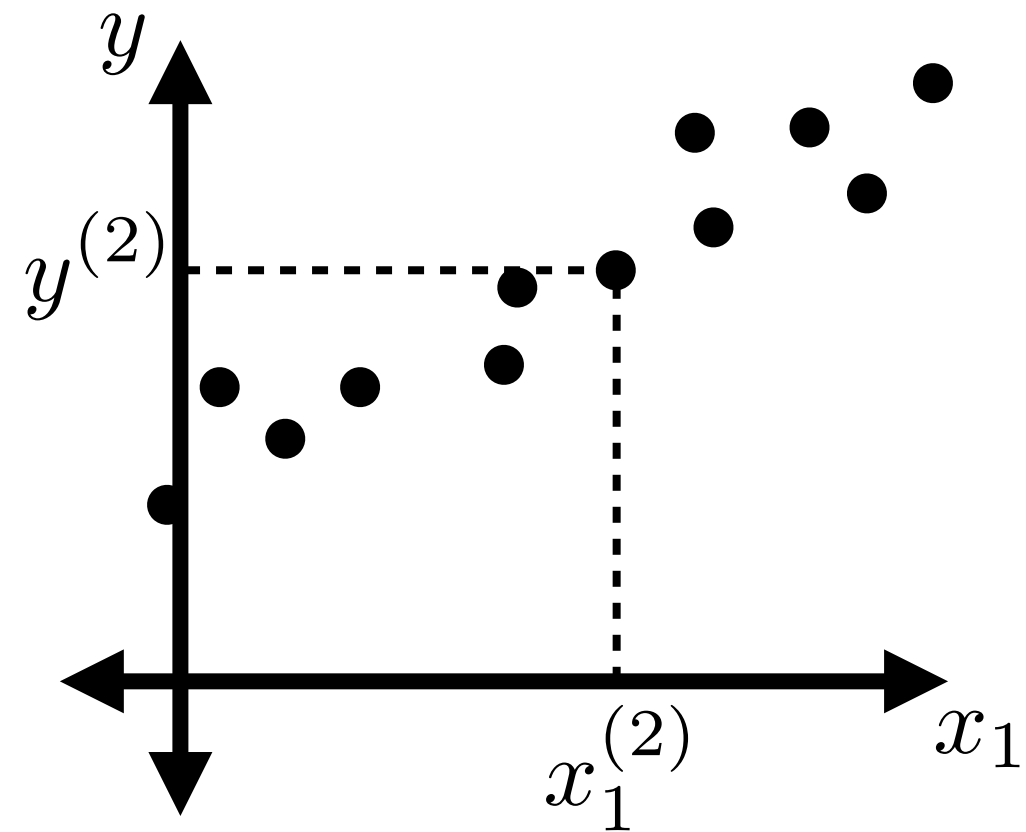
- Have data; have hypothesis class





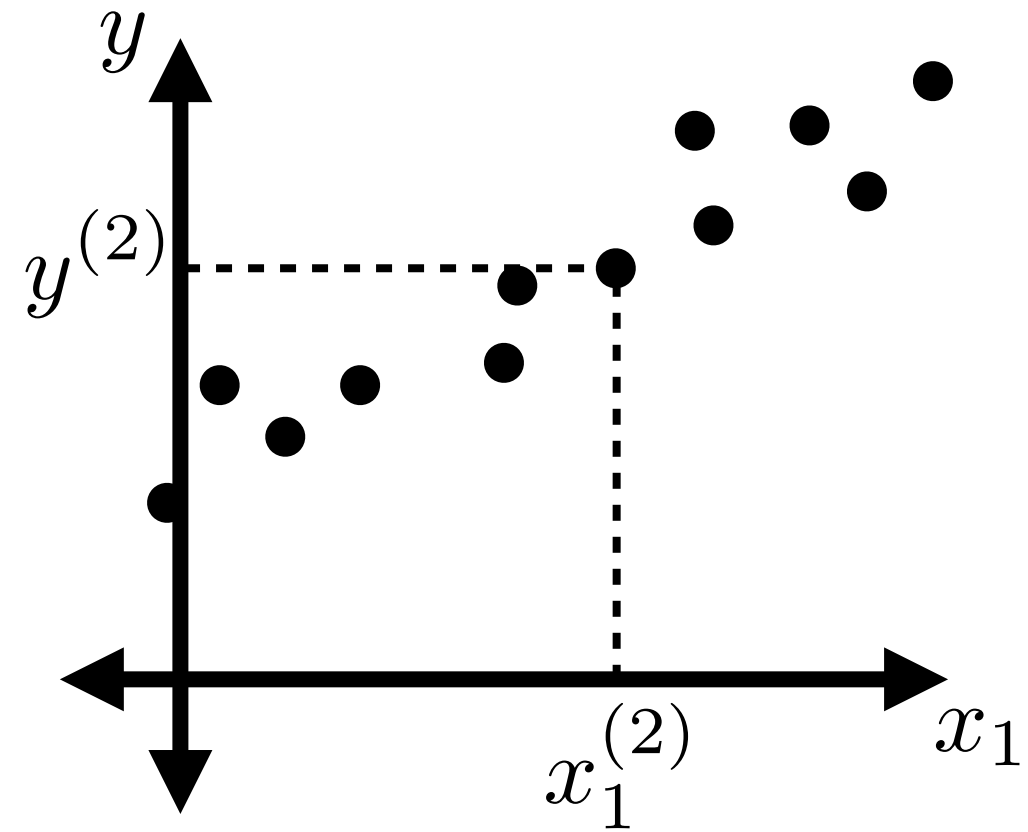
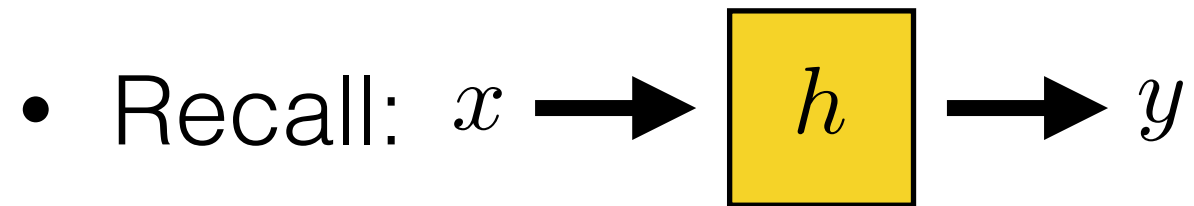
# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



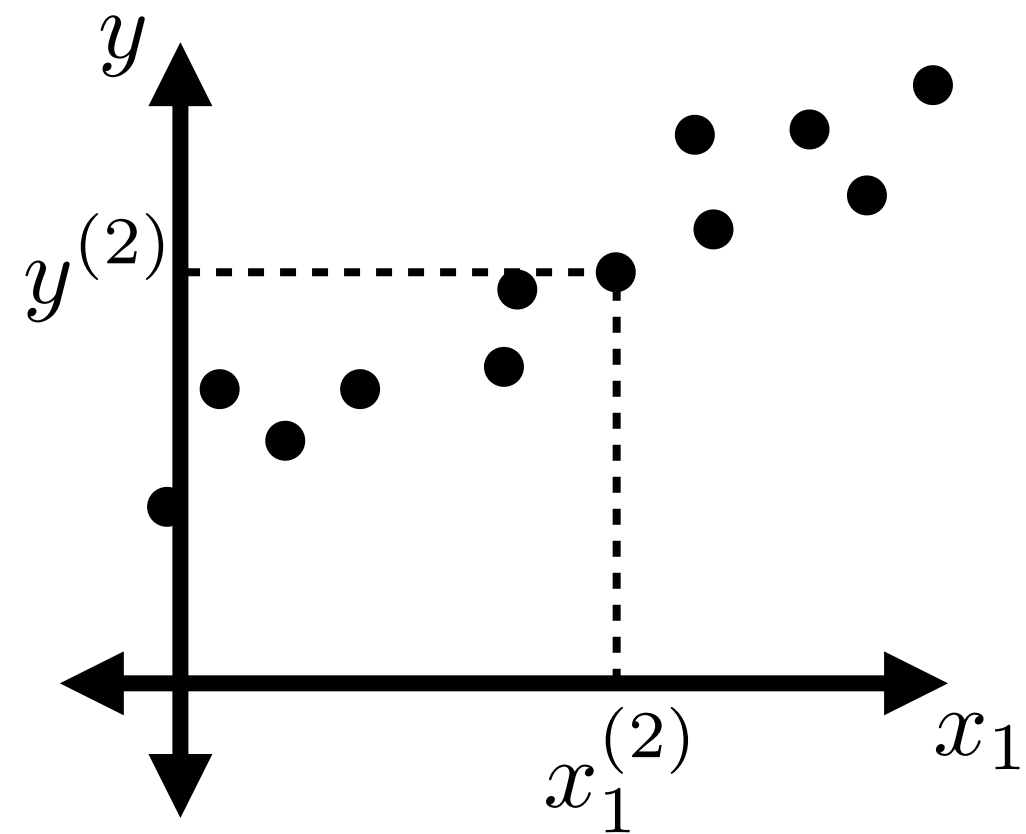
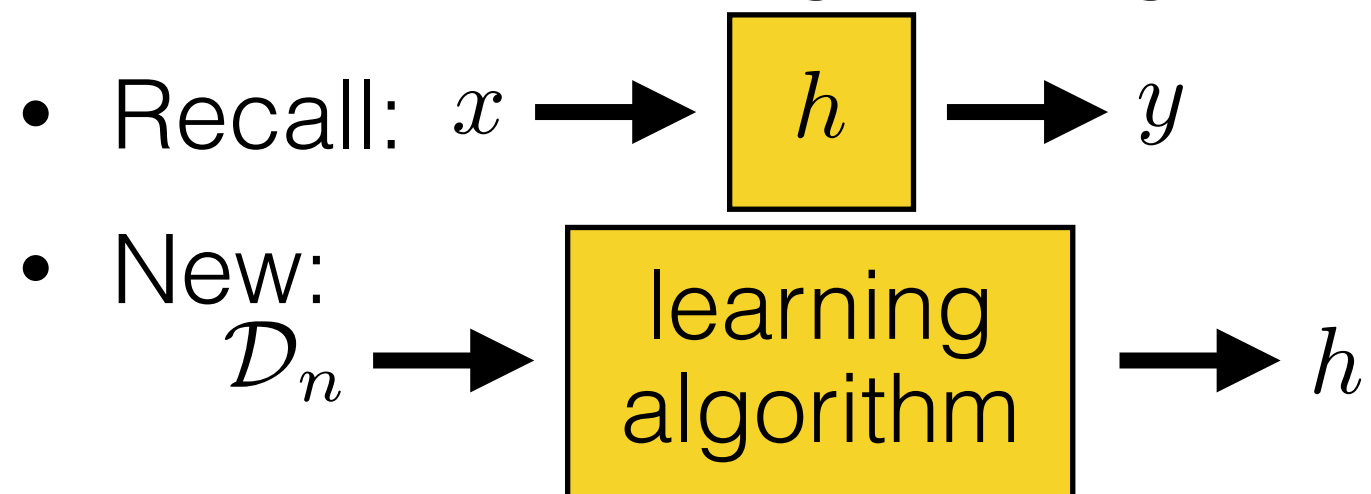
# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



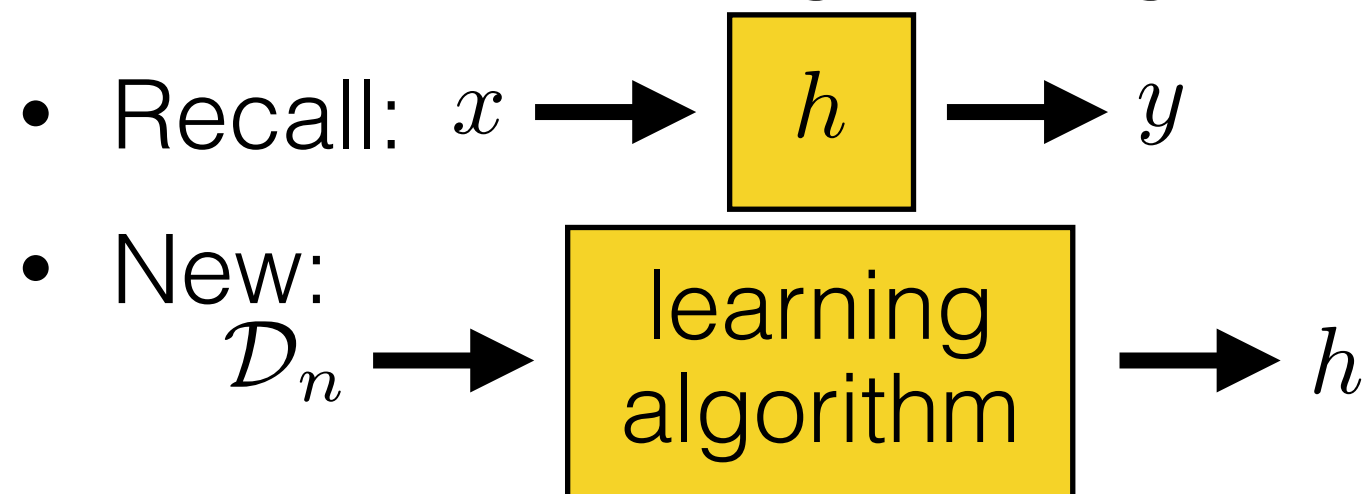
# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor

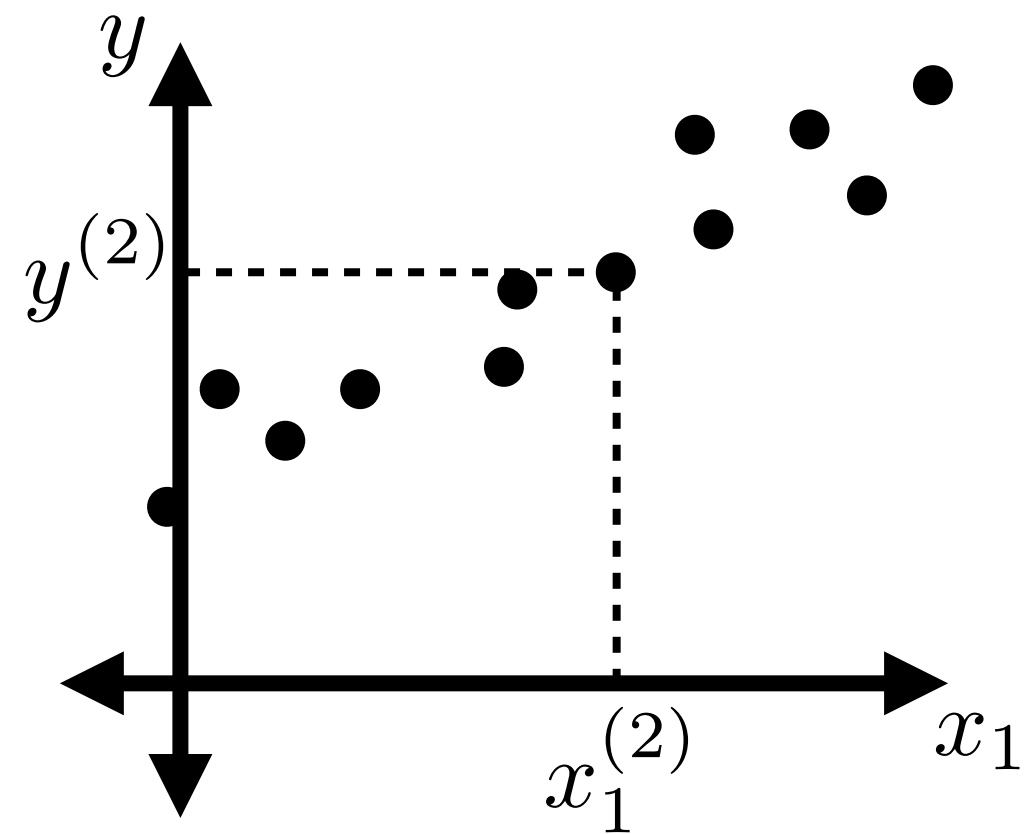


# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor

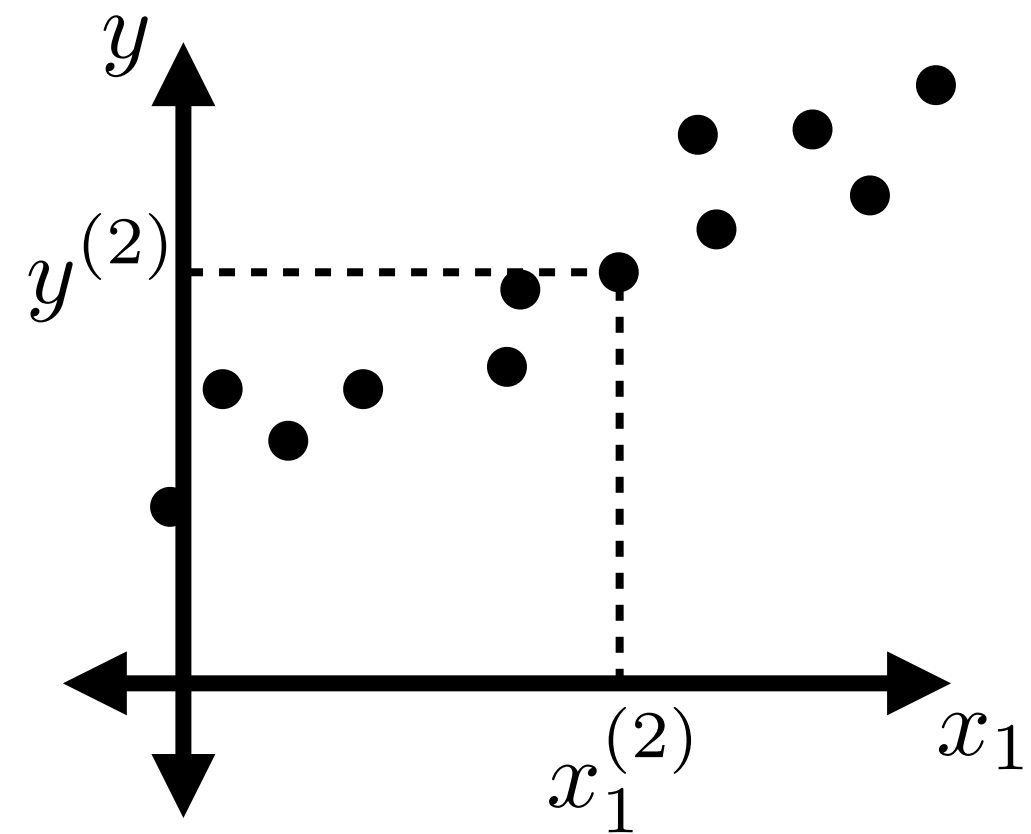
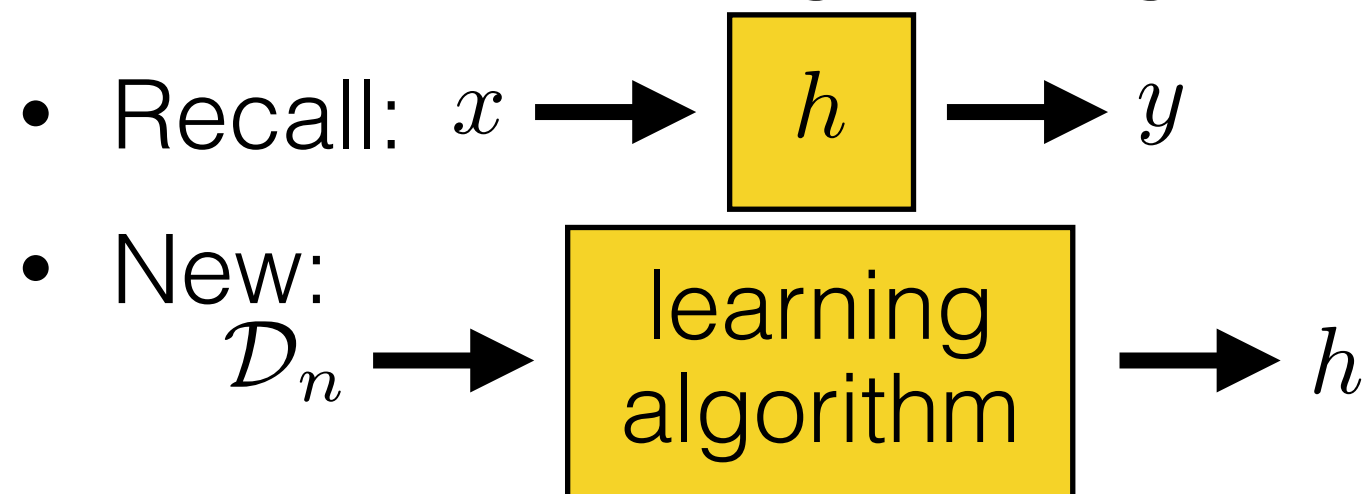


- Example:



# Learning a regressor

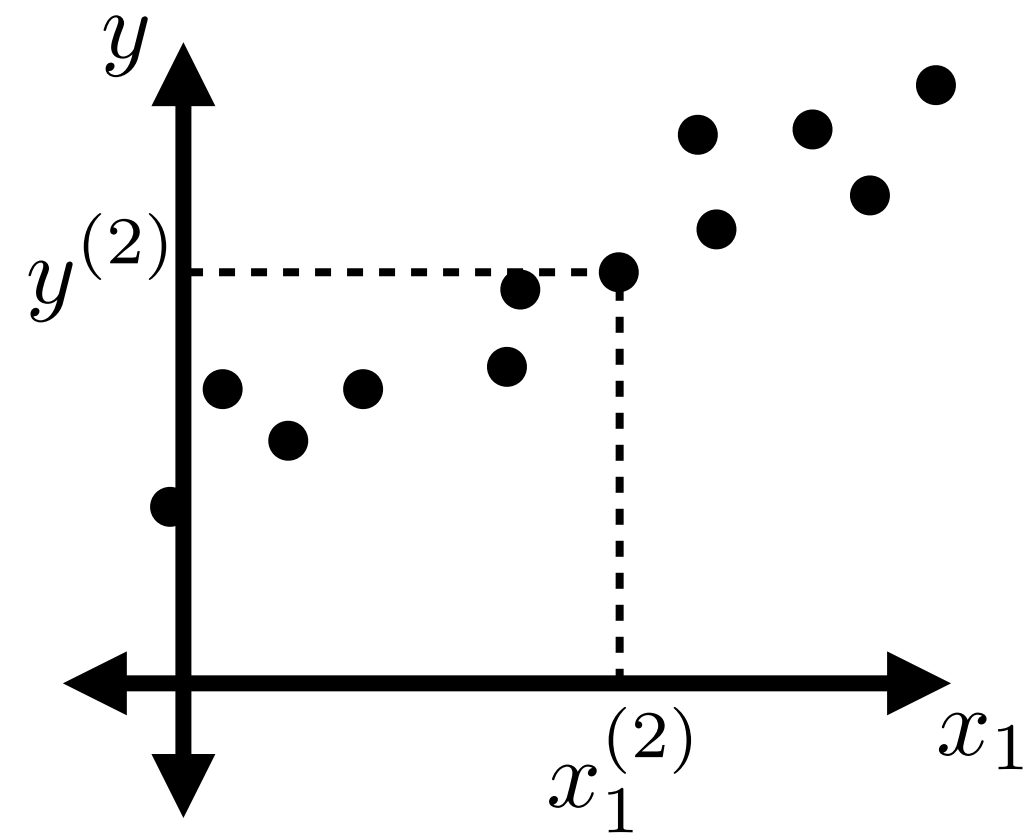
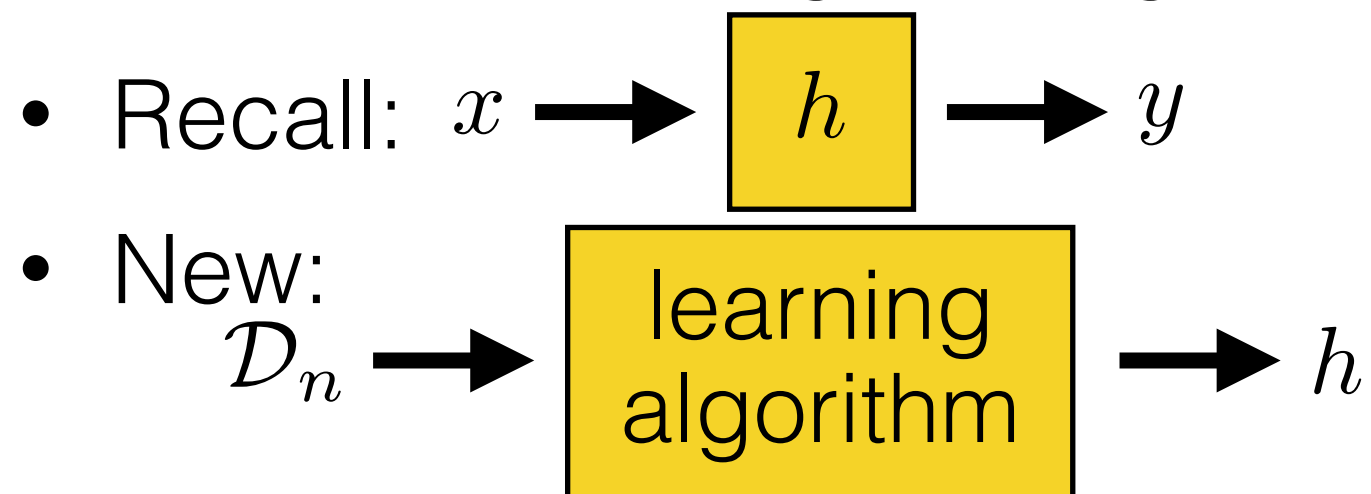
- Have data; have hypothesis class
- Want to choose a good regressor



- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



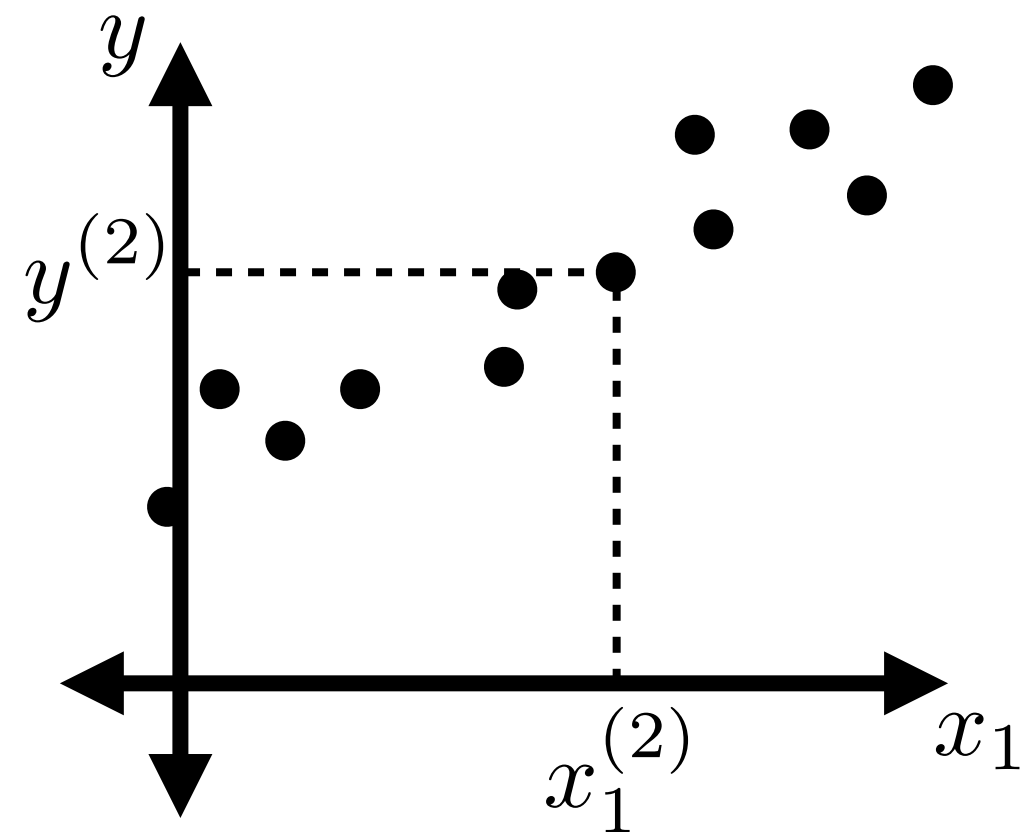
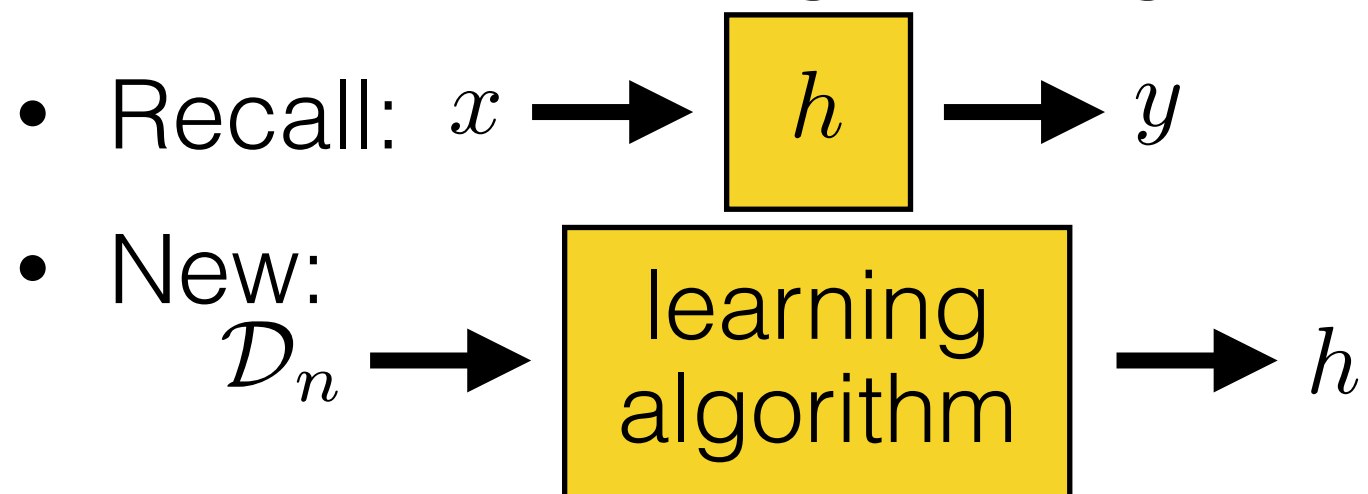
- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$



# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



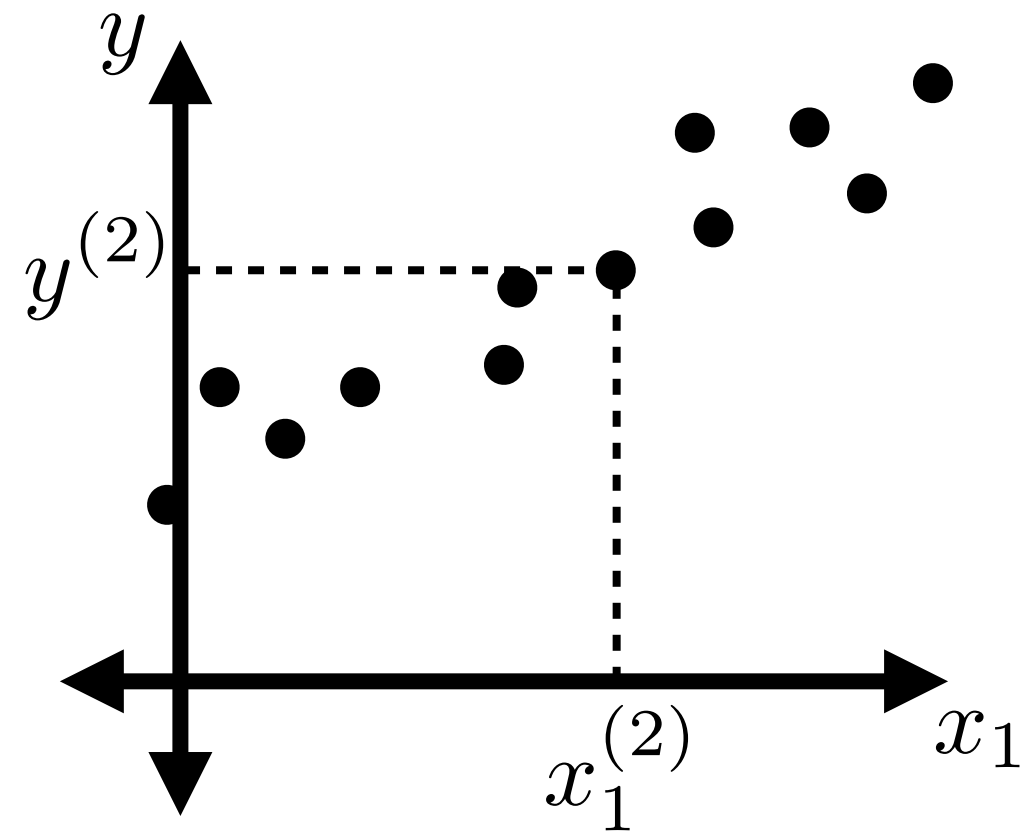
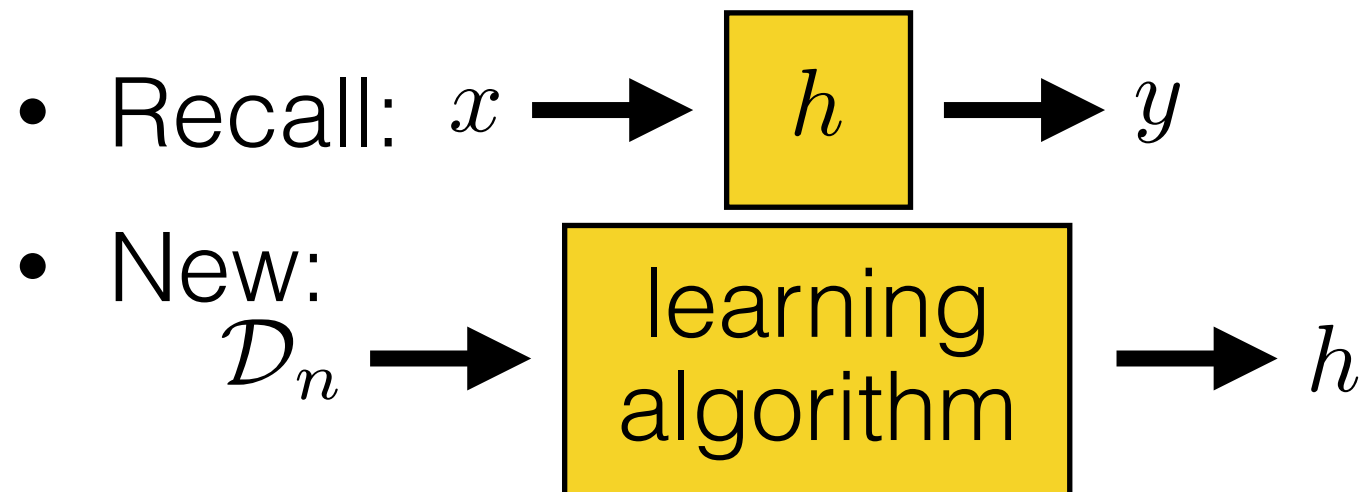
- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$



# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



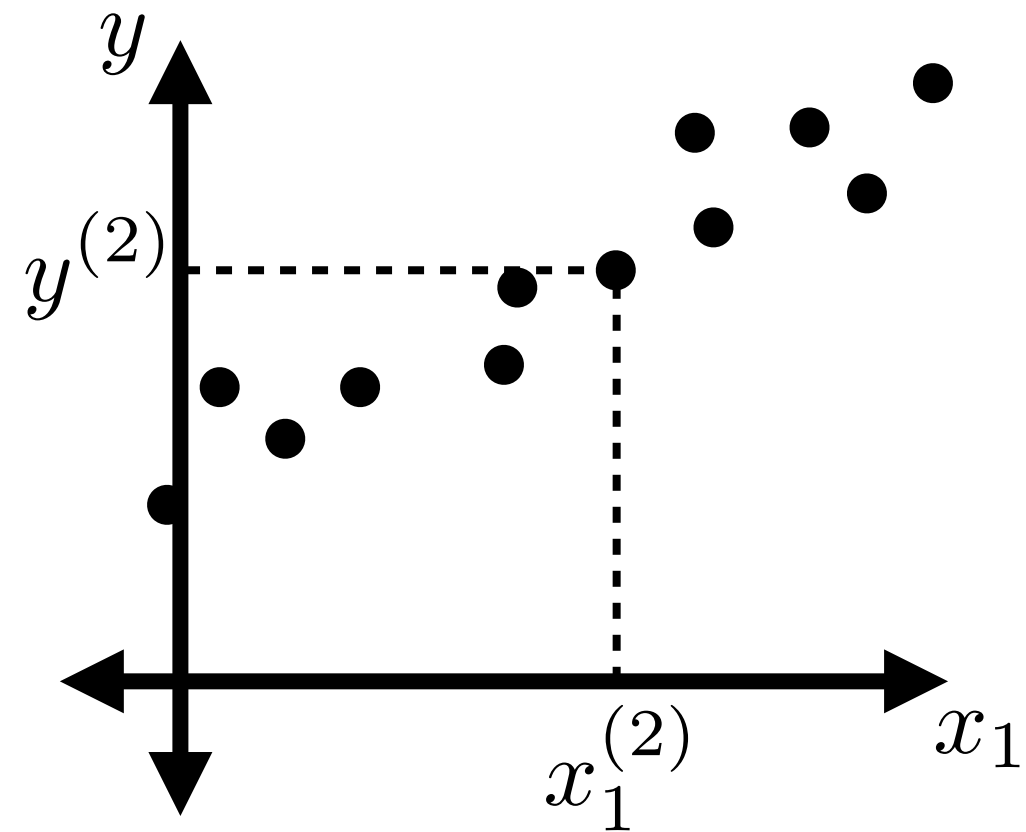
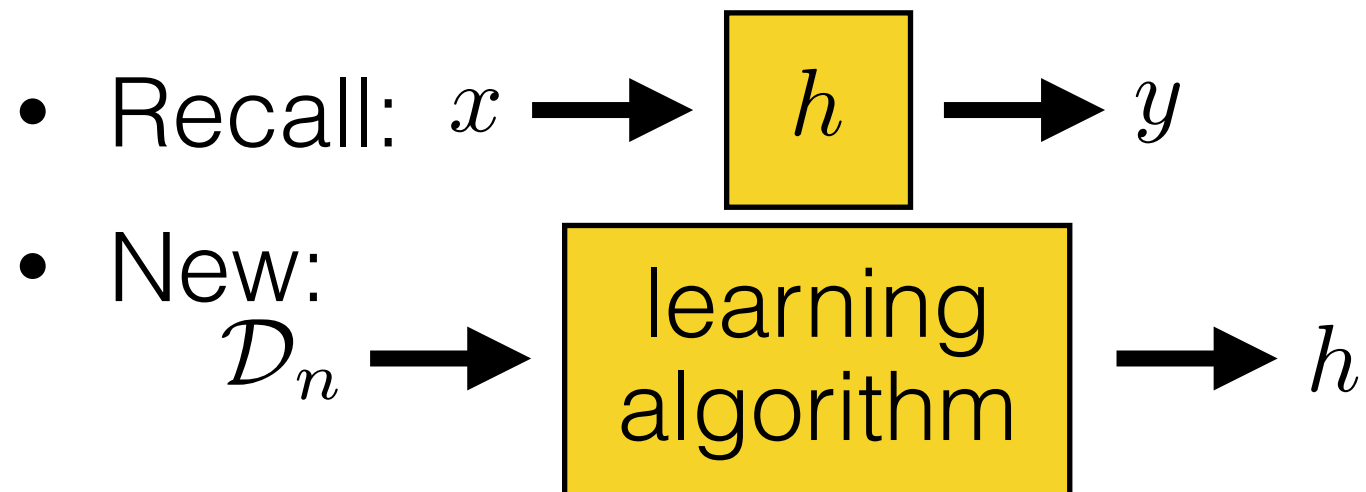
- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$

Ex\_learning\_alg(  $\mathcal{D}_n ; k$  )

# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

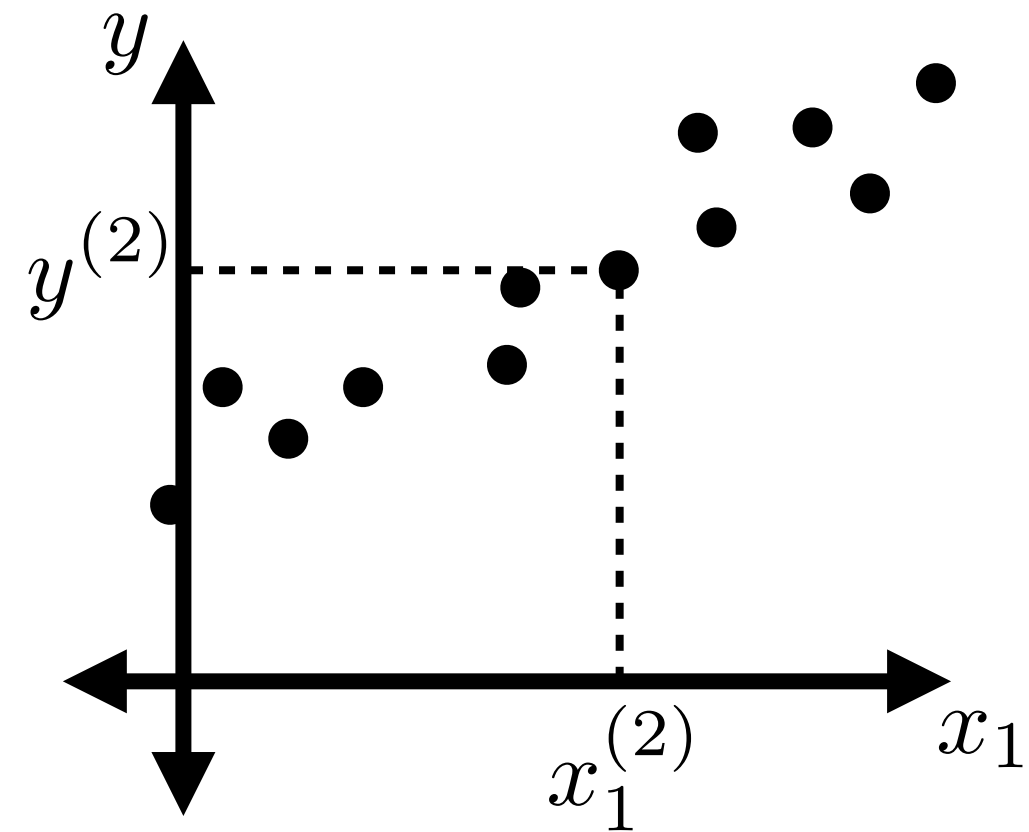
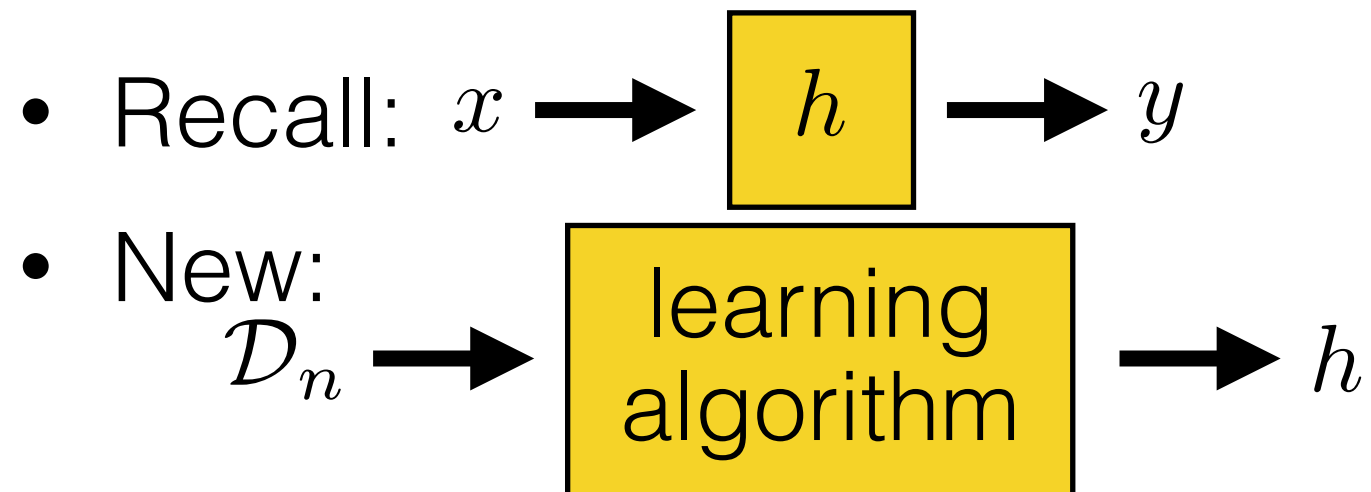
$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$

Ex\_learning\_alg( $\mathcal{D}_n; k$ )

Set  $j^* = \text{the } j \in \{1, \dots, k\} \text{ with lowest } \mathcal{E}_n(h^{(j)})$

# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$

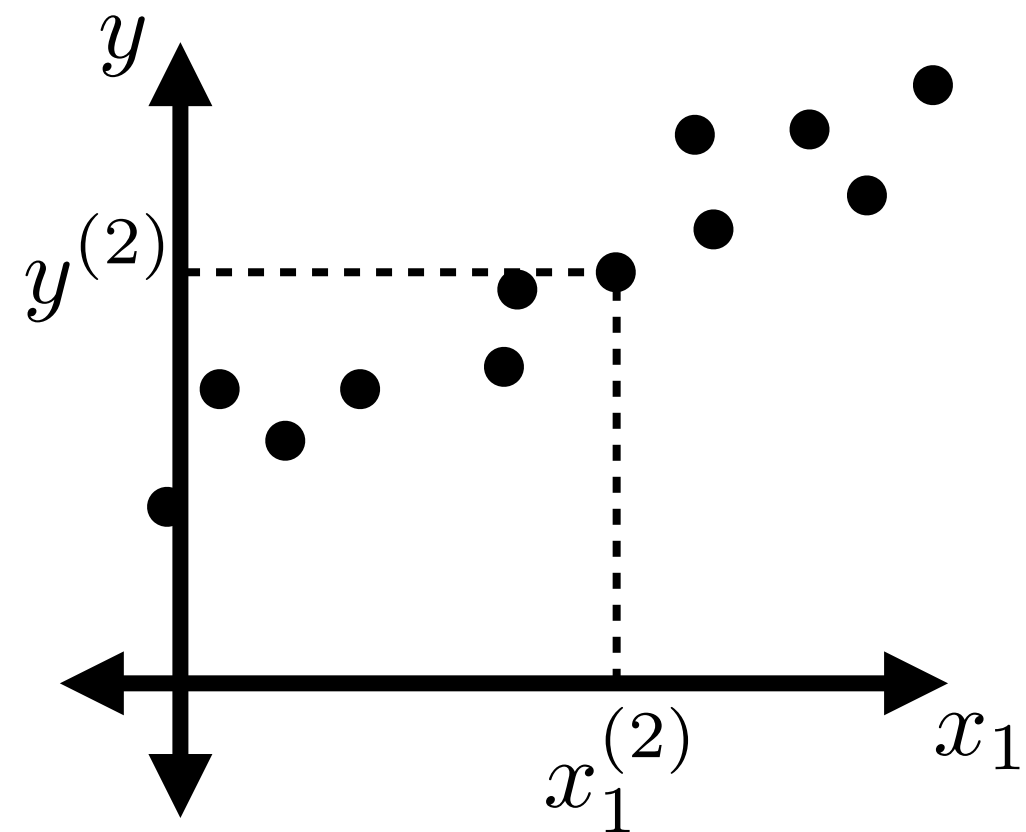
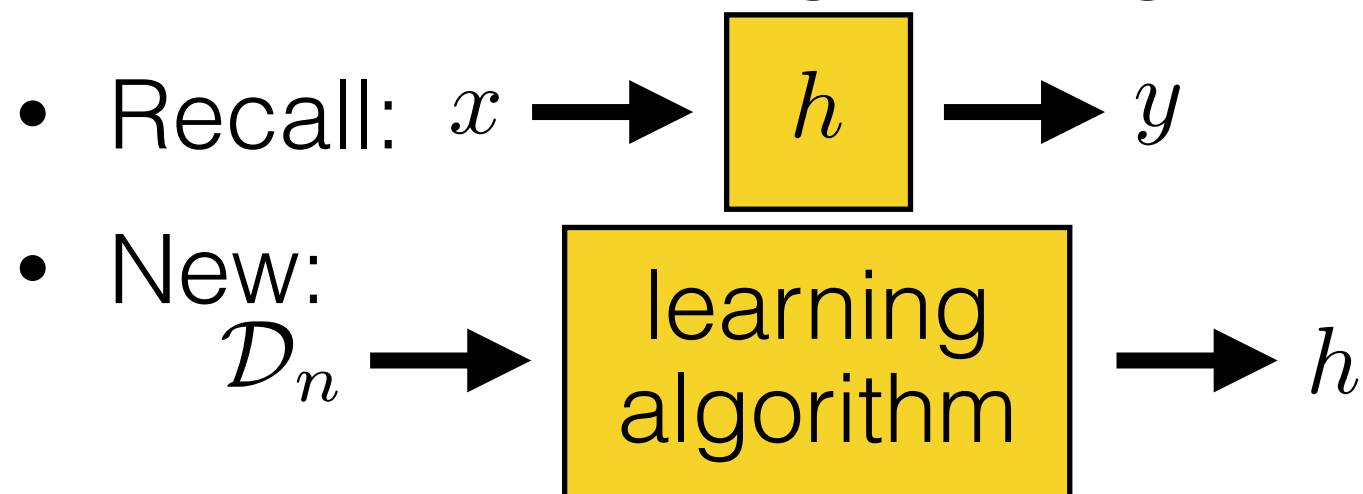
Ex\_learning\_alg(  $\mathcal{D}_n ; k$  )

Set  $j^* = \text{the } j \in \{1, \dots, k\} \text{ with lowest } \mathcal{E}_n(h^{(j)})$

Return  $h^{(j^*)}$

# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$

`Ex_learning_alg(  $\mathcal{D}_n$  ;  $k$  )`

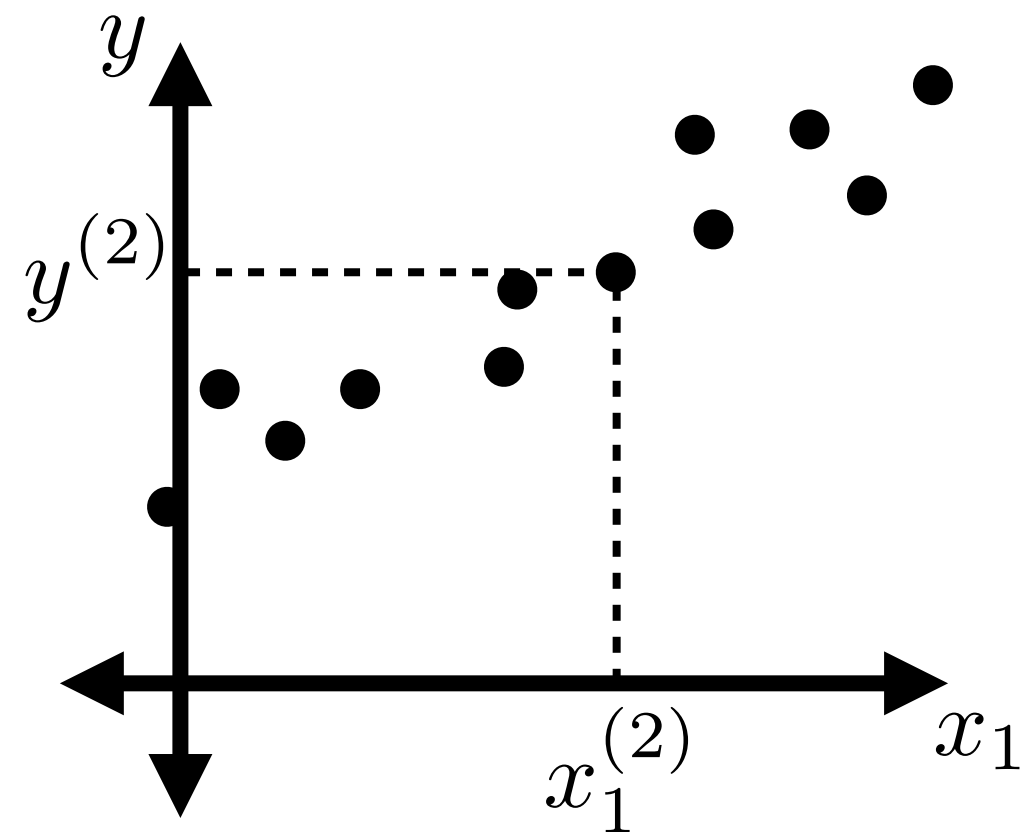
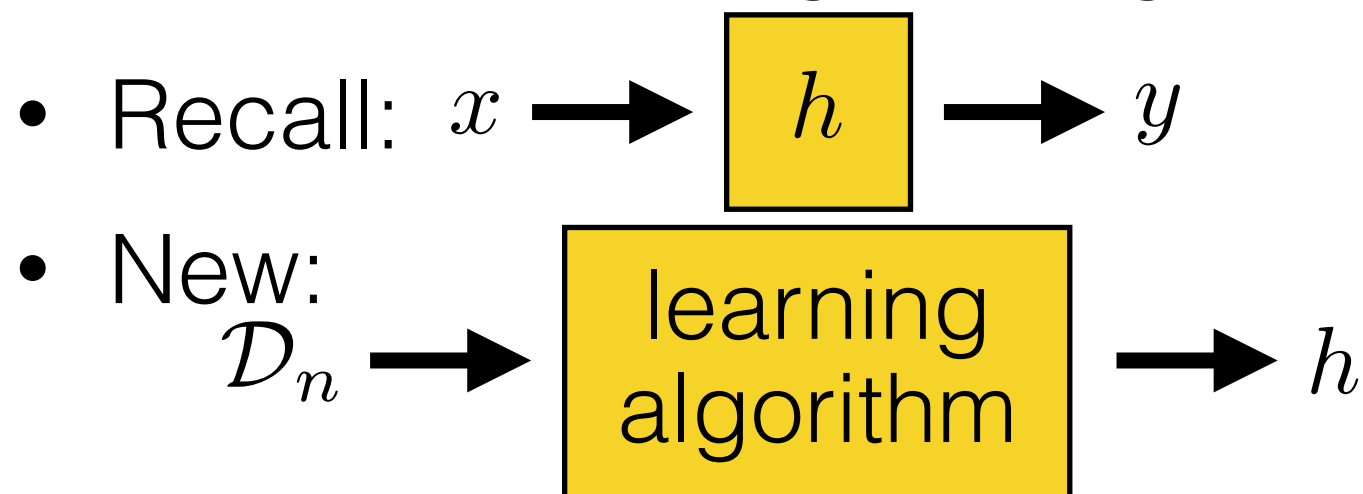
Set  $j^* = \text{the } j \in \{1, \dots, k\} \text{ with lowest } \mathcal{E}_n(h^{(j)})$

Return  $h^{(j^*)}$

hyperparameter

# Learning a regressor

- Have data; have hypothesis class
- Want to choose a good regressor



- Example:
  - Suppose someone already generated 1 trillion hypotheses, e.g. at random, indexed by  $j$ :

$$h^{(j)}(x) = h(x; \theta^{(j)}, \theta_0^{(j)})$$

hyperparameter

`Ex_learning_alg(  $\mathcal{D}_n$ ;  $k$  )`

Set  $j^* = \text{the } j \in \{1, \dots, k\} \text{ with lowest } \mathcal{E}_n(h^{(j)})$

Return  $h^{(j^*)}$

- How does training error of `Ex_learning_alg( $\mathcal{D}_n$ ;1)` compare to the training error of `Ex_learning_alg( $\mathcal{D}_n$ ;2)`?

# Linear regression: Another way



# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error: 
$$\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss

$$\frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$\frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$\frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$$



# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$\frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} =$  
$$\begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} =$  
$$\begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$   $\tilde{Y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

$n \times d$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$   $\tilde{Y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

*nx d*      *nx 1*



# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$   $\tilde{Y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2 = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$   $\tilde{Y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

# Linear regression: Another way

- How about we just consider all hypotheses in our class and choose the one with lowest training error?
  - We'll see: not typically straightforward
  - But for linear regression with square loss: can do it!

- Recall: training error:  $\mathcal{E}_n(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^{(i)}), y^{(i)})$

- Training error: square loss, linear regr., extra "1" feature

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} - y^{(i)})^2 = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$$

Define  $\tilde{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$   $\tilde{Y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

$n \times d$        $n \times 1$

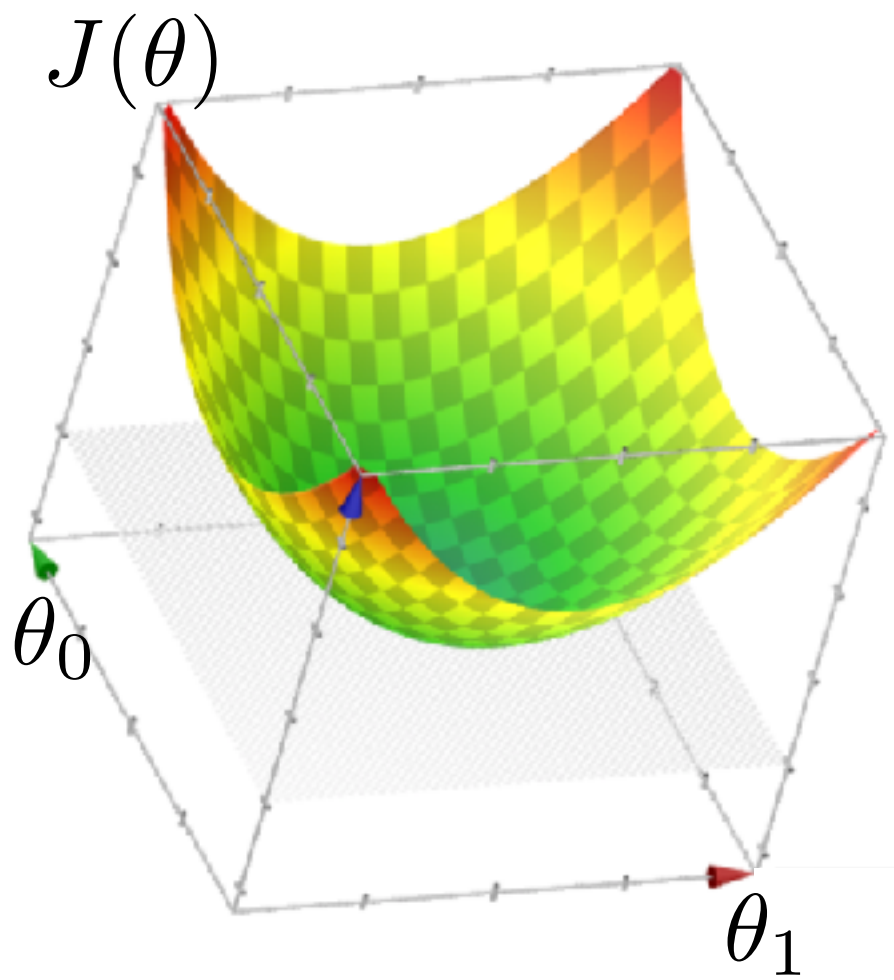
- Goal: minimize  $J(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$

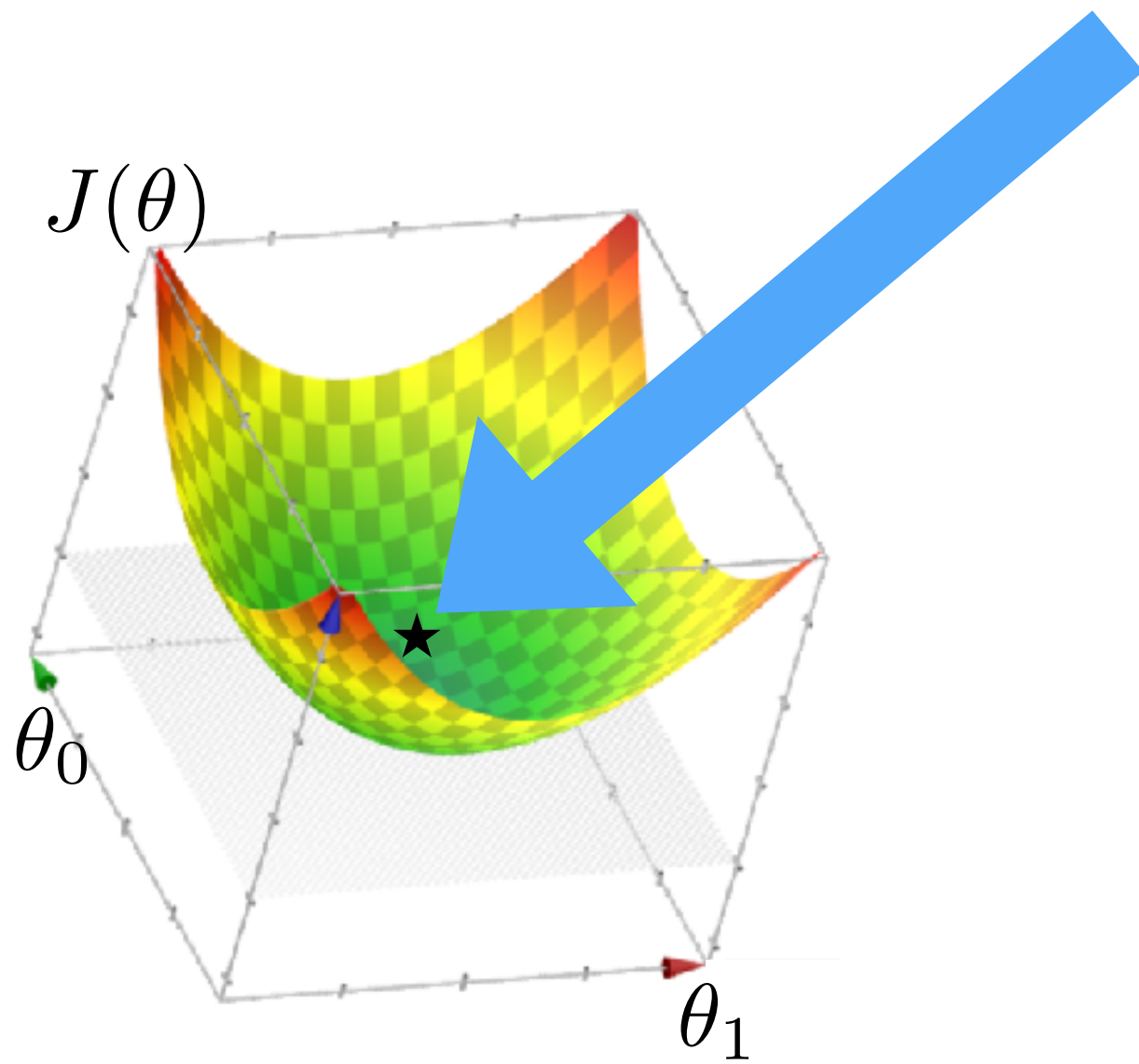
# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$



# Linear regression: A Direct Solution

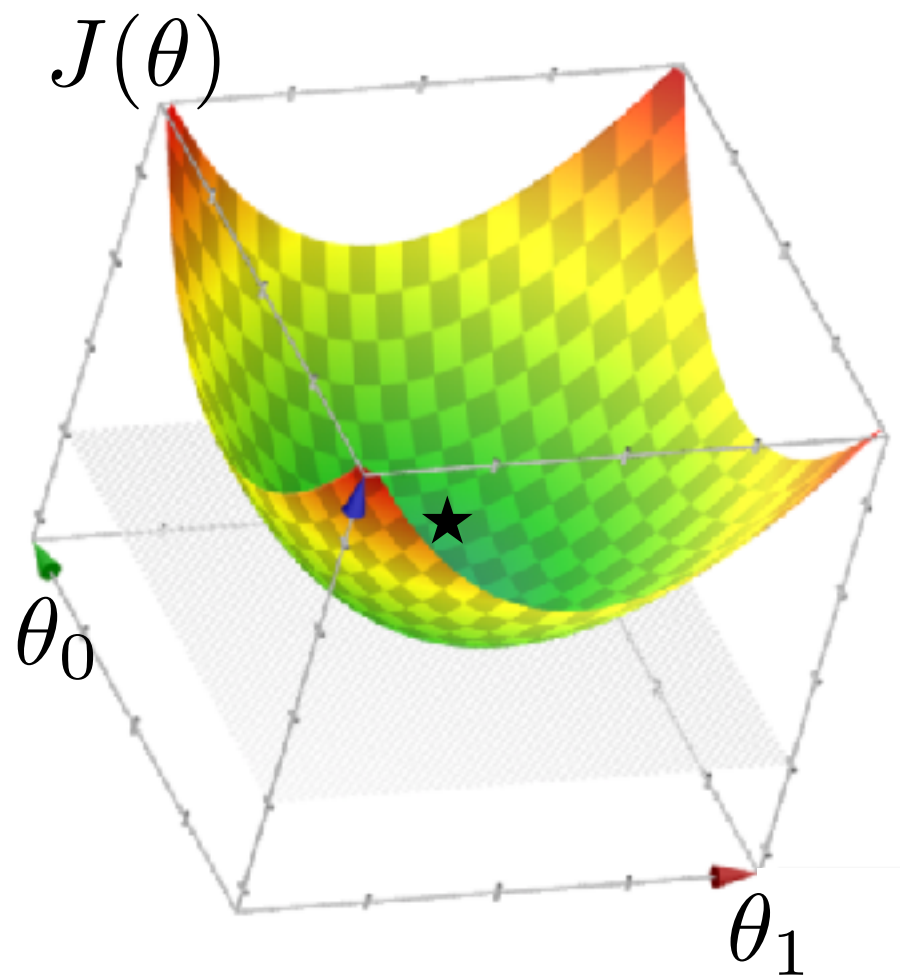
- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$





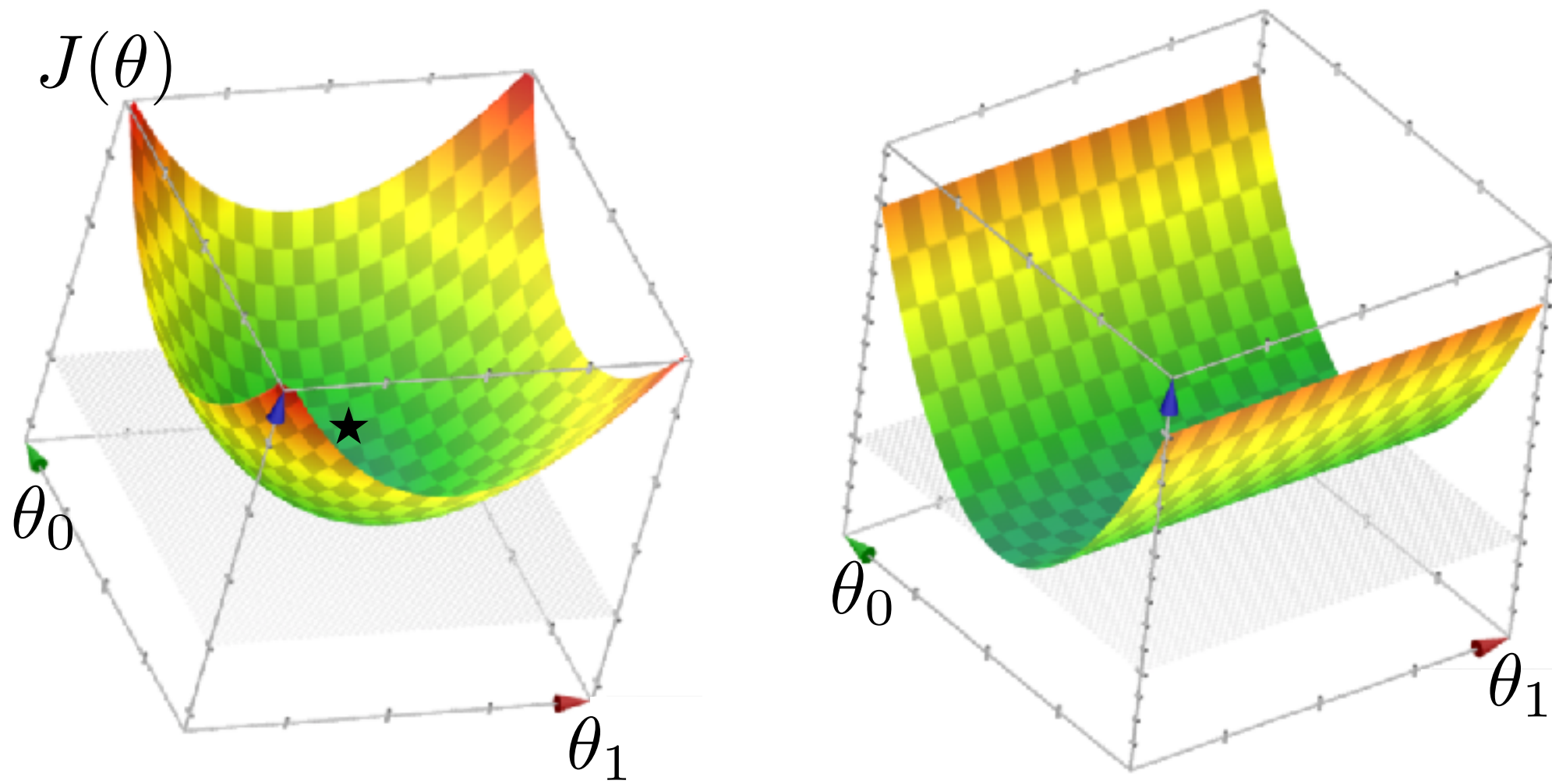
# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$



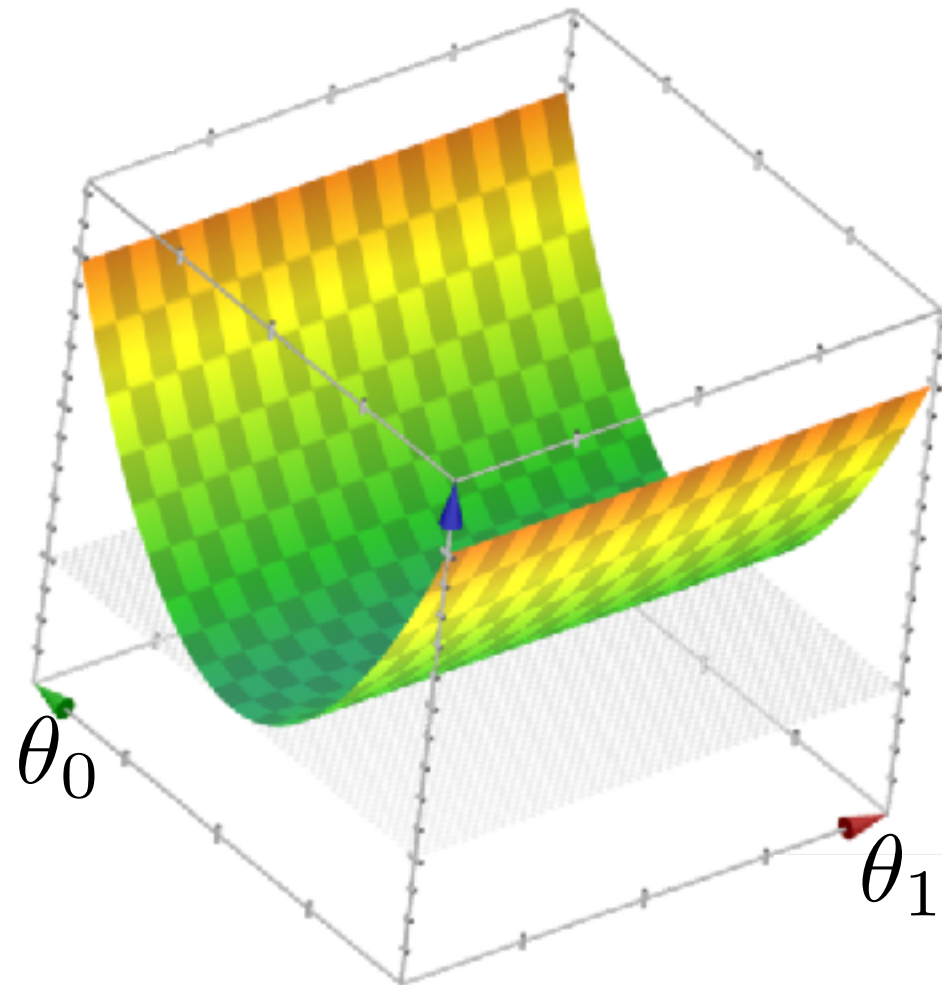
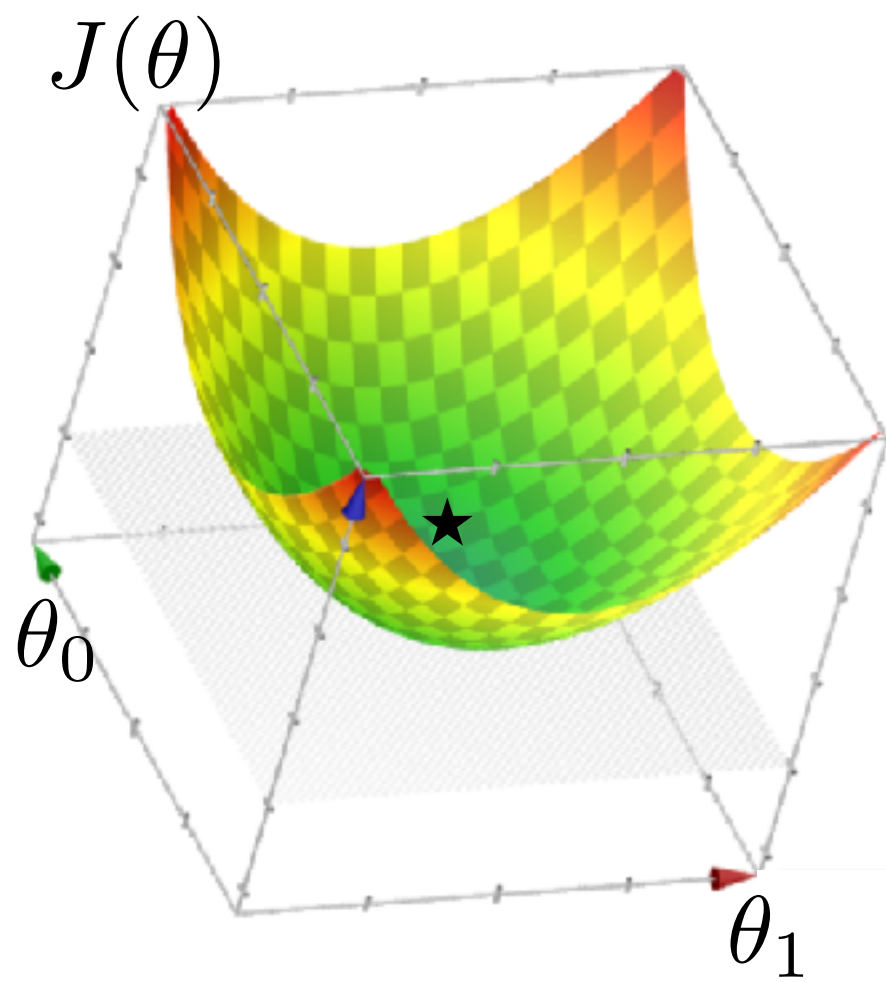
# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta)$

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta)$

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0$  set



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} \text{set } 0$

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

# Linear regression: A Direct Solution

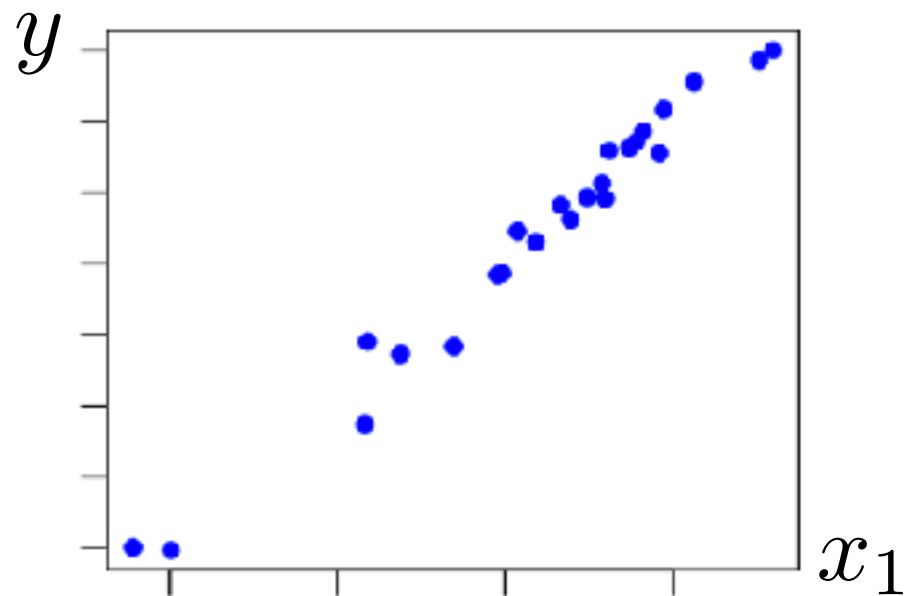
- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

Exercise:  
check  $n, d=1$

# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

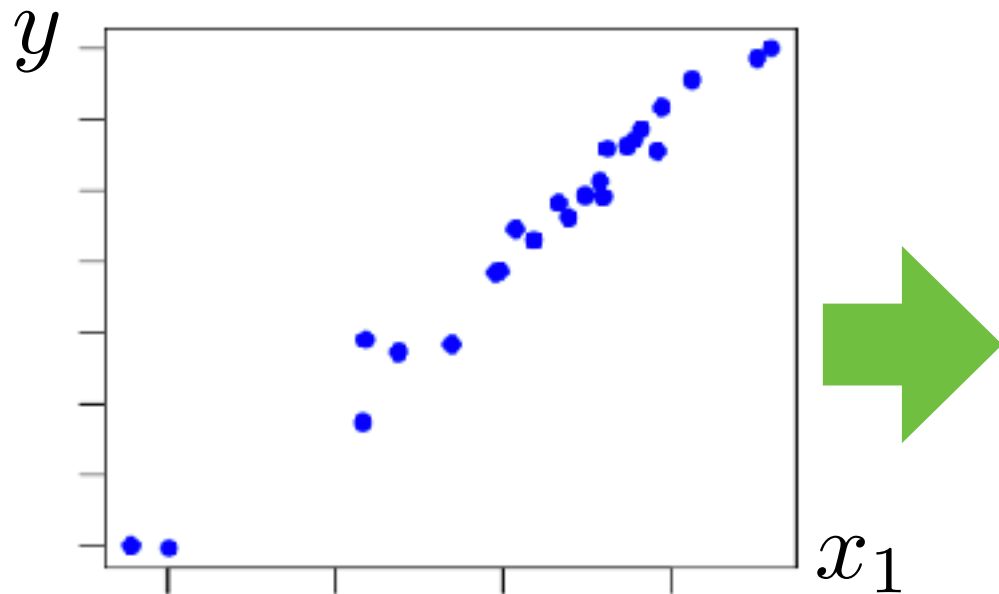
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top(\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

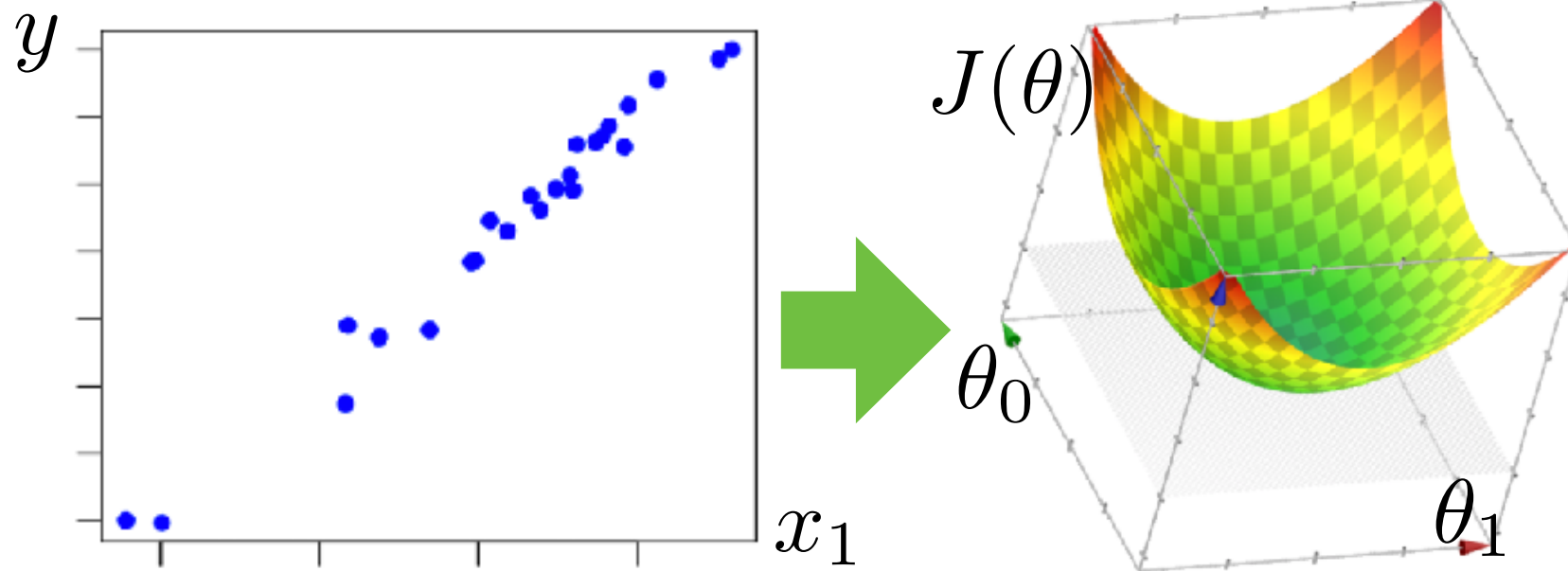
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

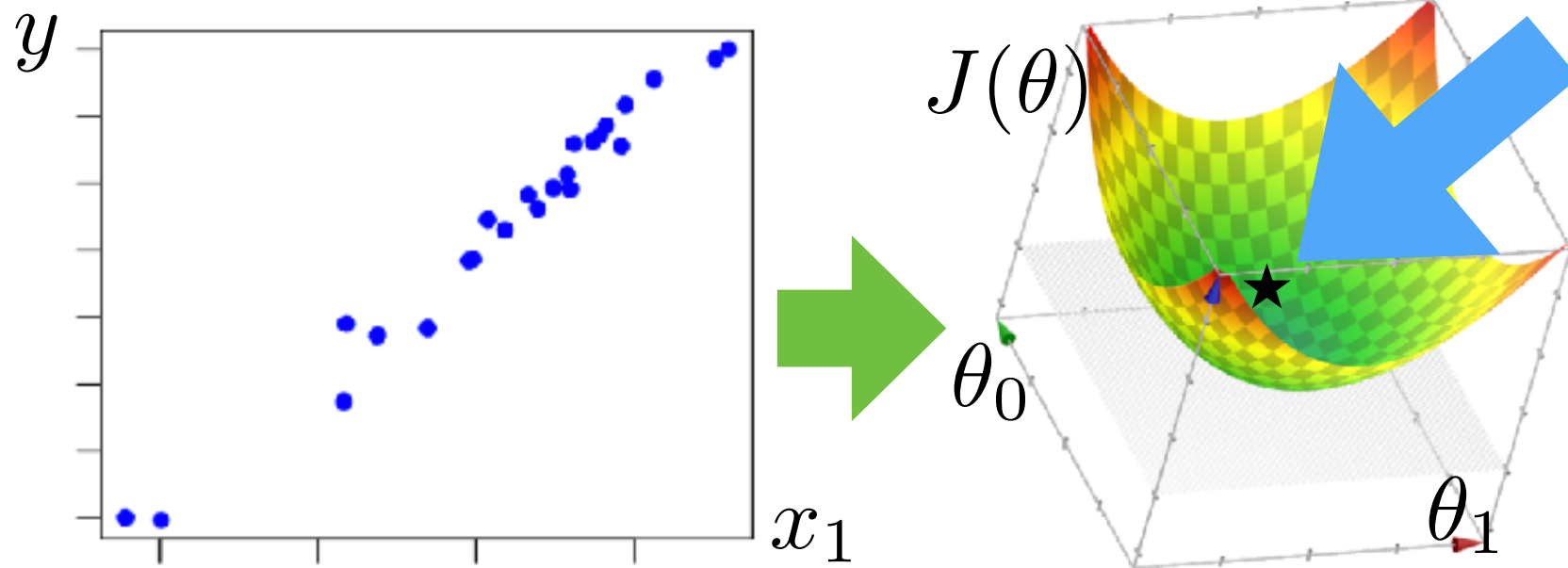
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top(\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

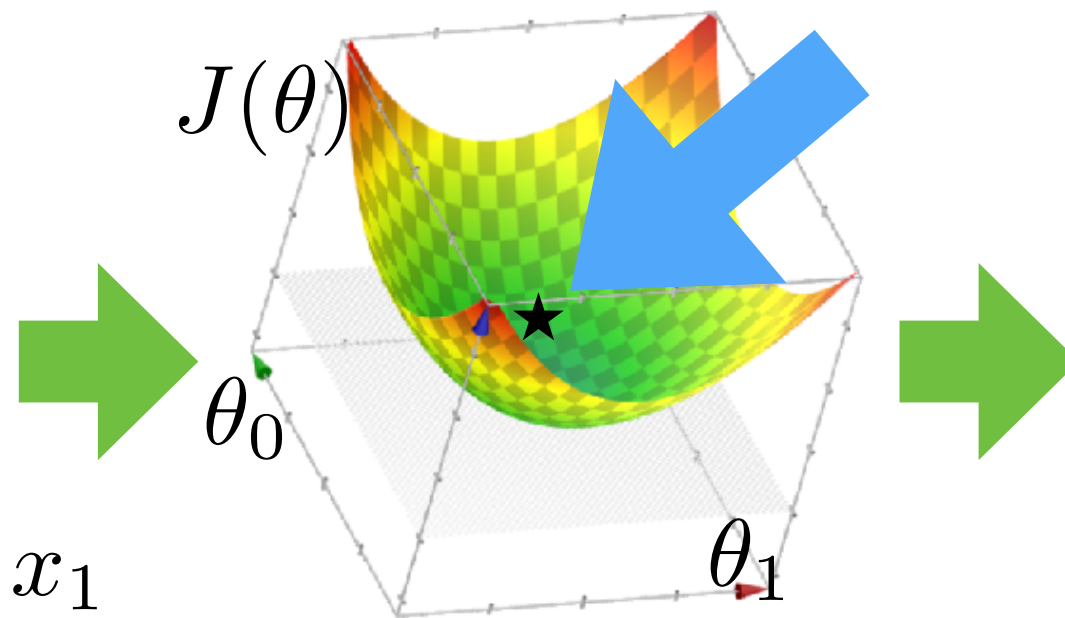
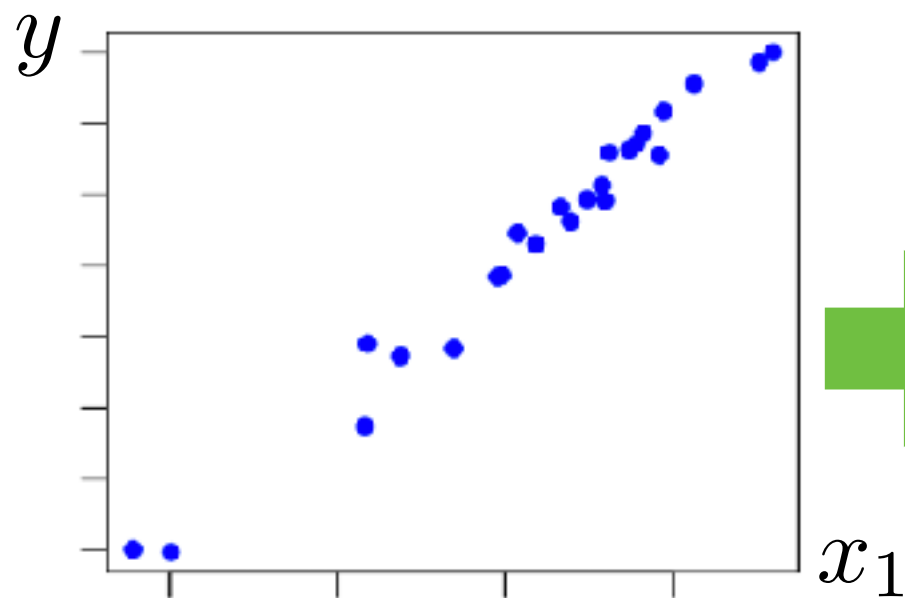
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top(\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

Exercise:  
check  $n, d=1$

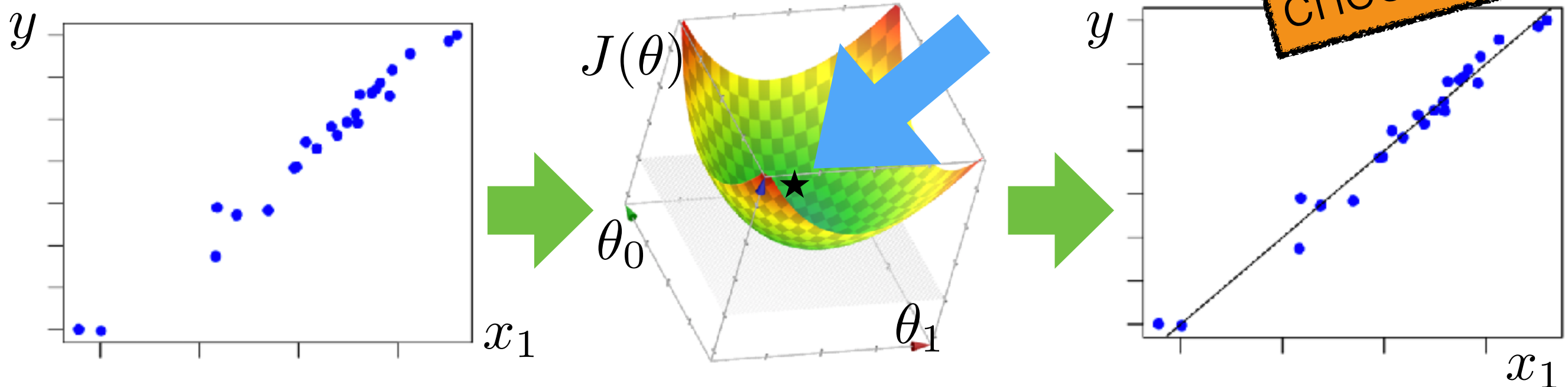




# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1}}{=} 0 \stackrel{\text{set}}{\Rightarrow} \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

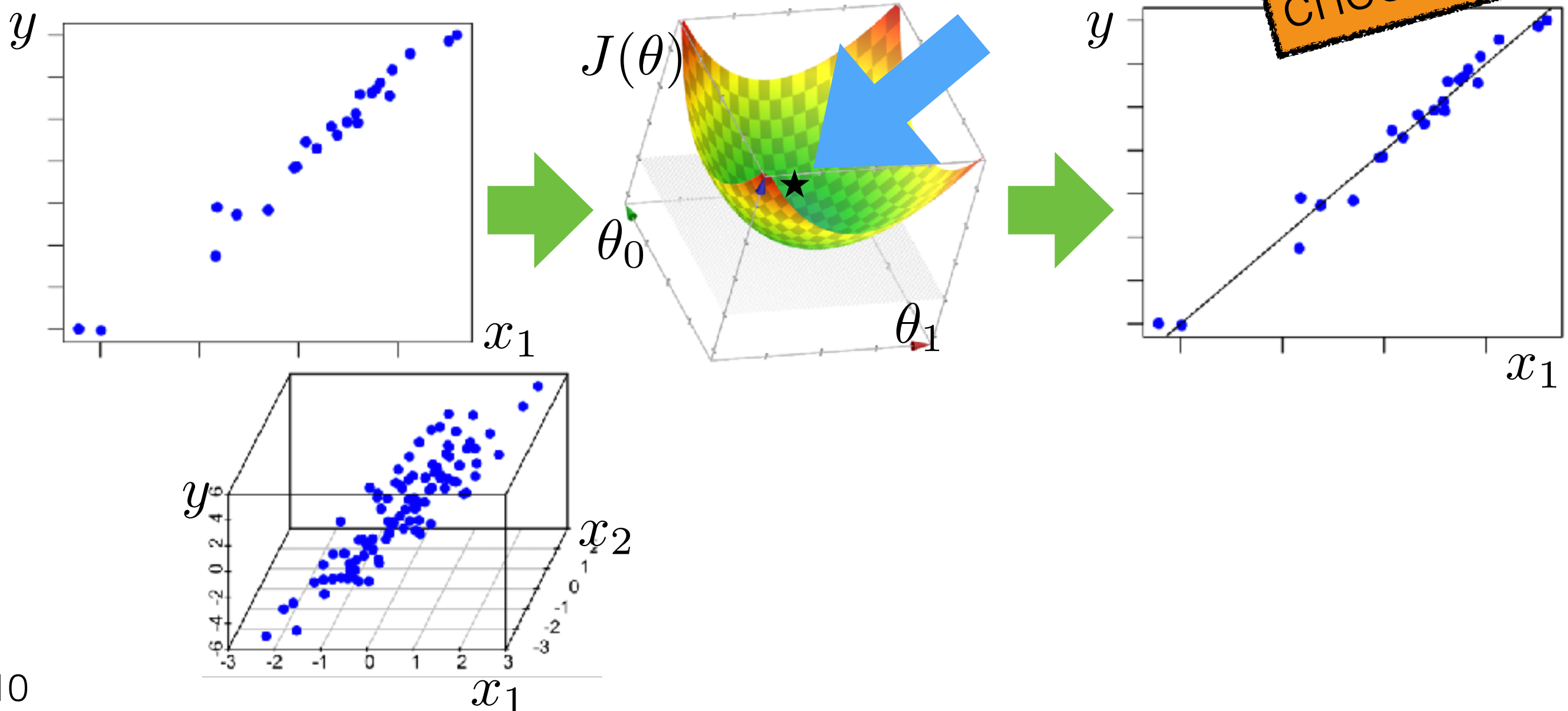
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top(\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1 set}}{=} 0 \Rightarrow \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

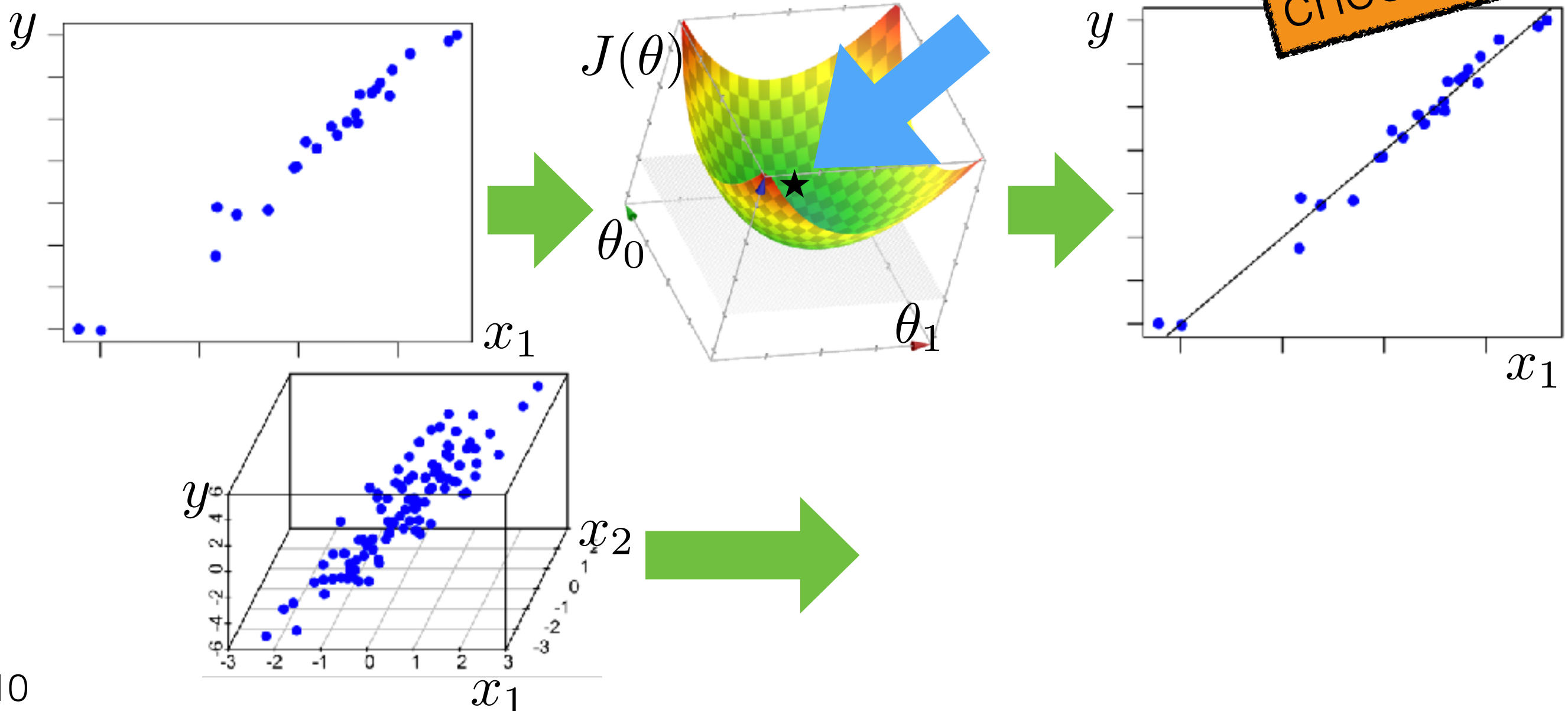
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top(\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1 set}}{=} 0 \Rightarrow \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

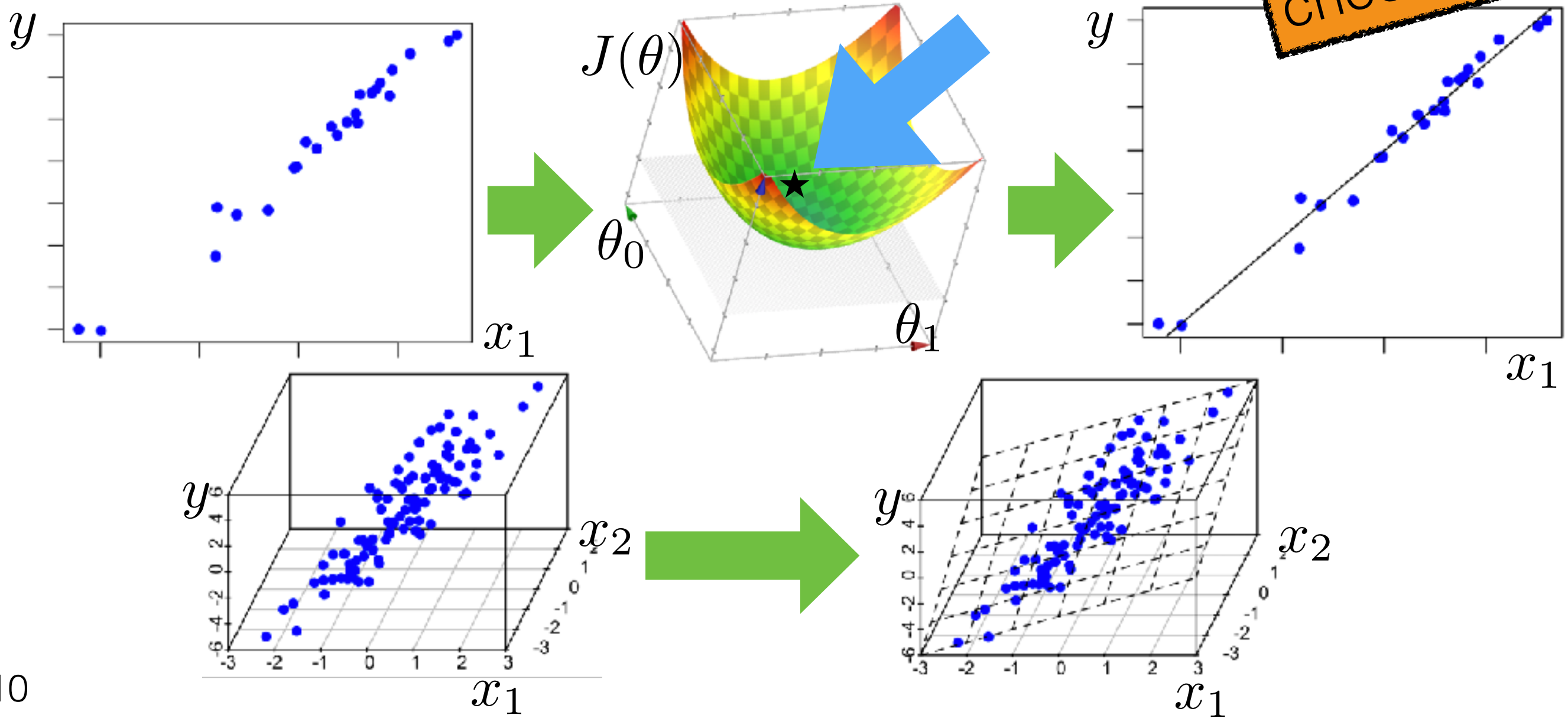
Exercise:  
check  $n, d=1$



# Linear regression: A Direct Solution

- Goal: minimize  $J(\theta) = \frac{1}{n}(\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y})$
- Uniquely minimized at a point if gradient at that point is zero and function “curves up” [see linear algebra]
- Gradient  $\nabla_{\theta} J(\theta) \stackrel{\text{dx1 set}}{=} 0 \Rightarrow \theta = (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y}$

Exercise:  
check  $n, d=1$



# What can go wrong in practice?

# What can go wrong in practice?

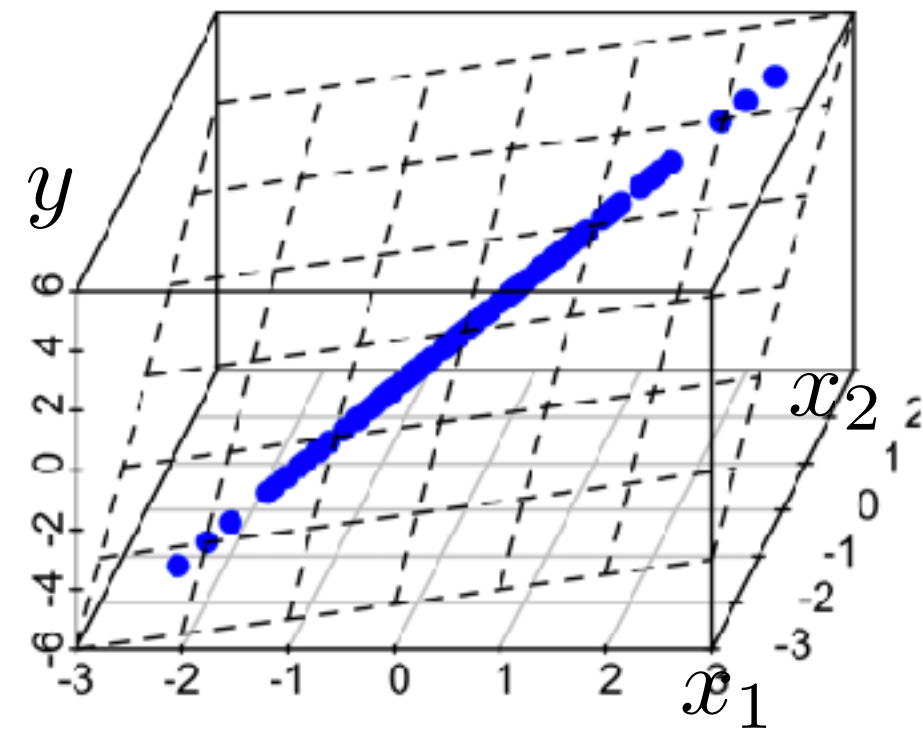
- Does the linear regr. objective always “curve up”? No!

# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane

# What can go wrong in practice?

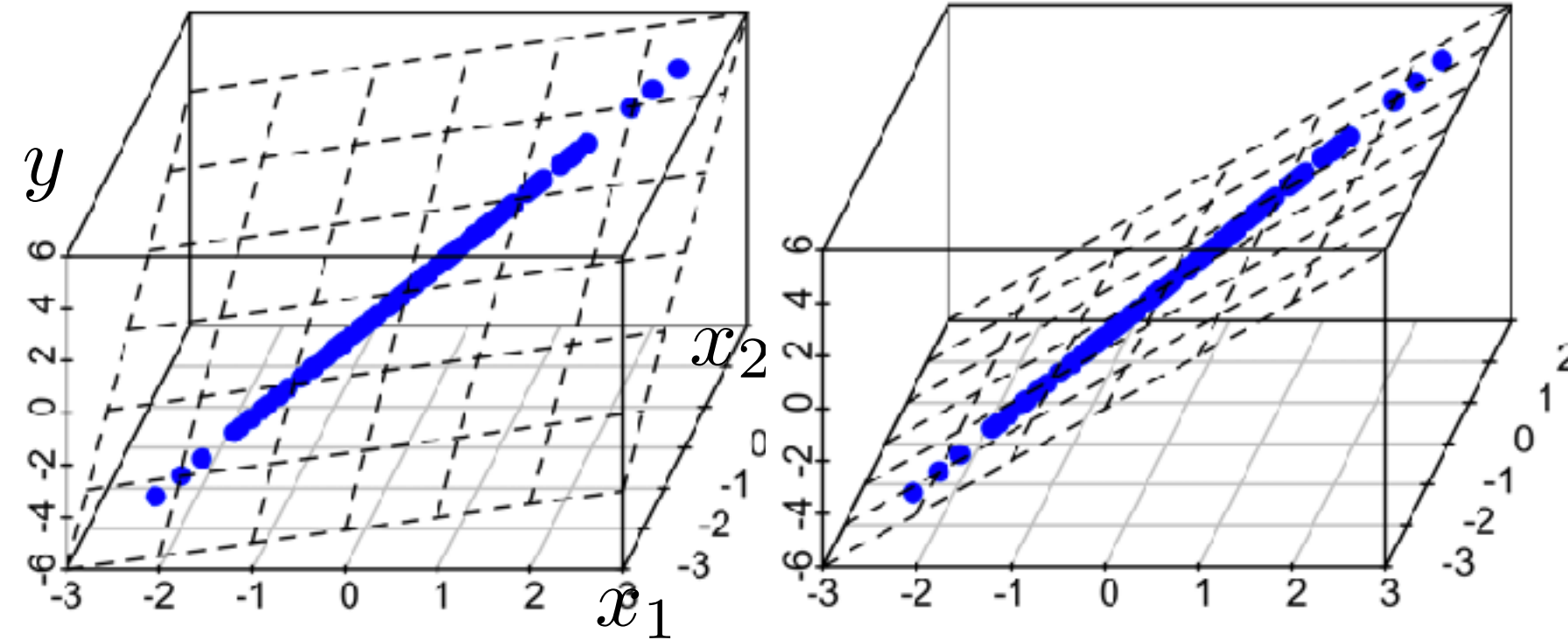
- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane





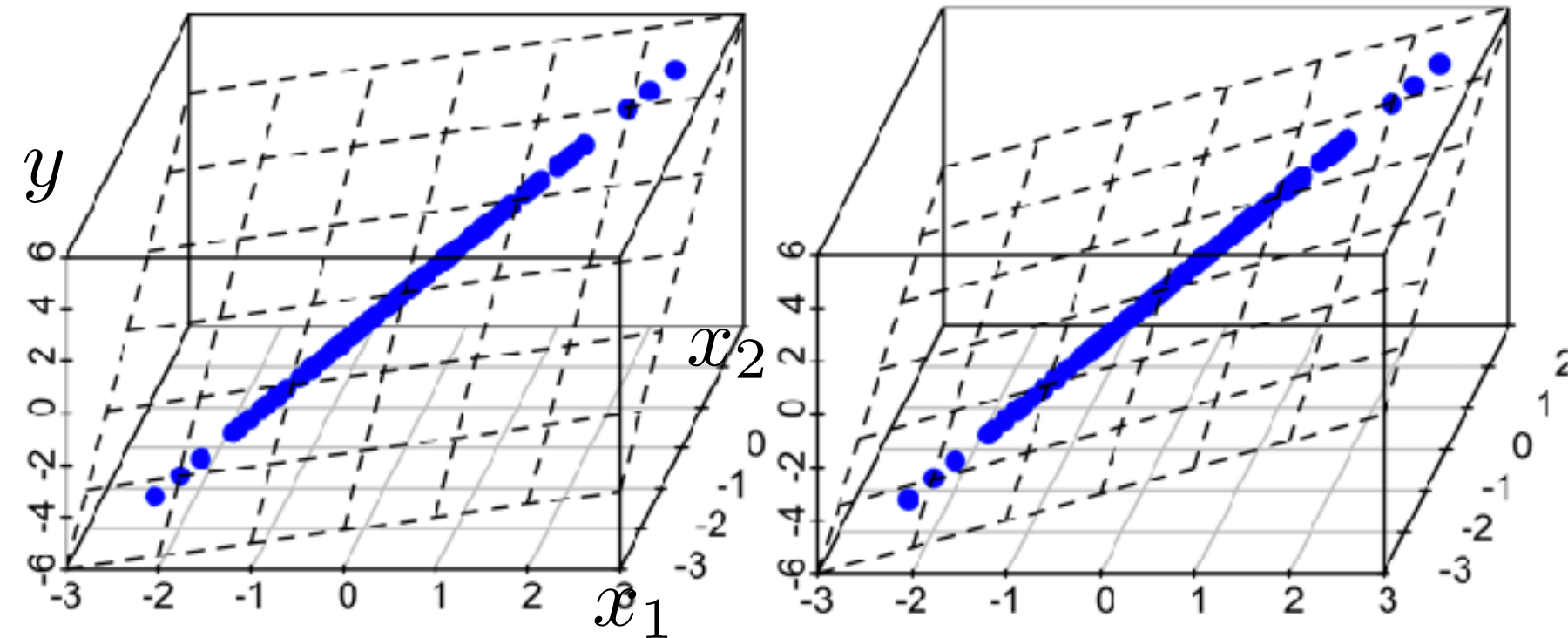
# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane



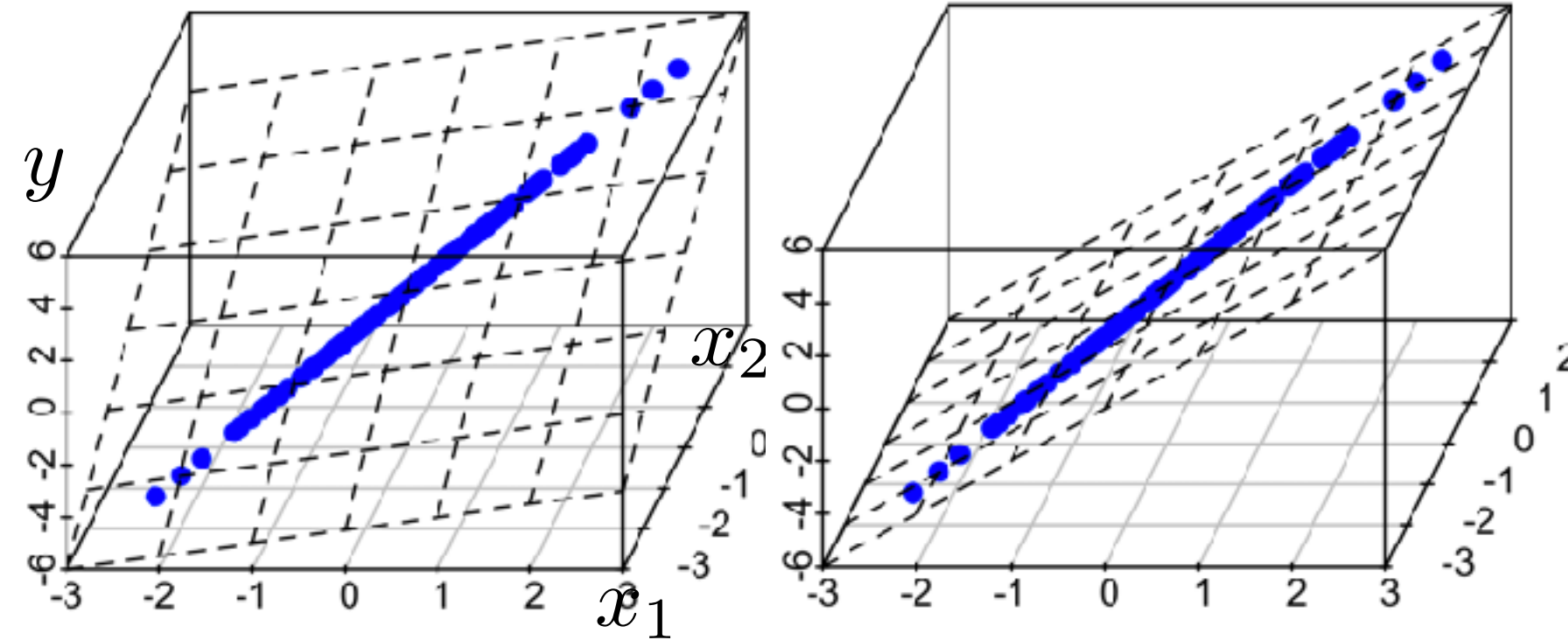
# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane



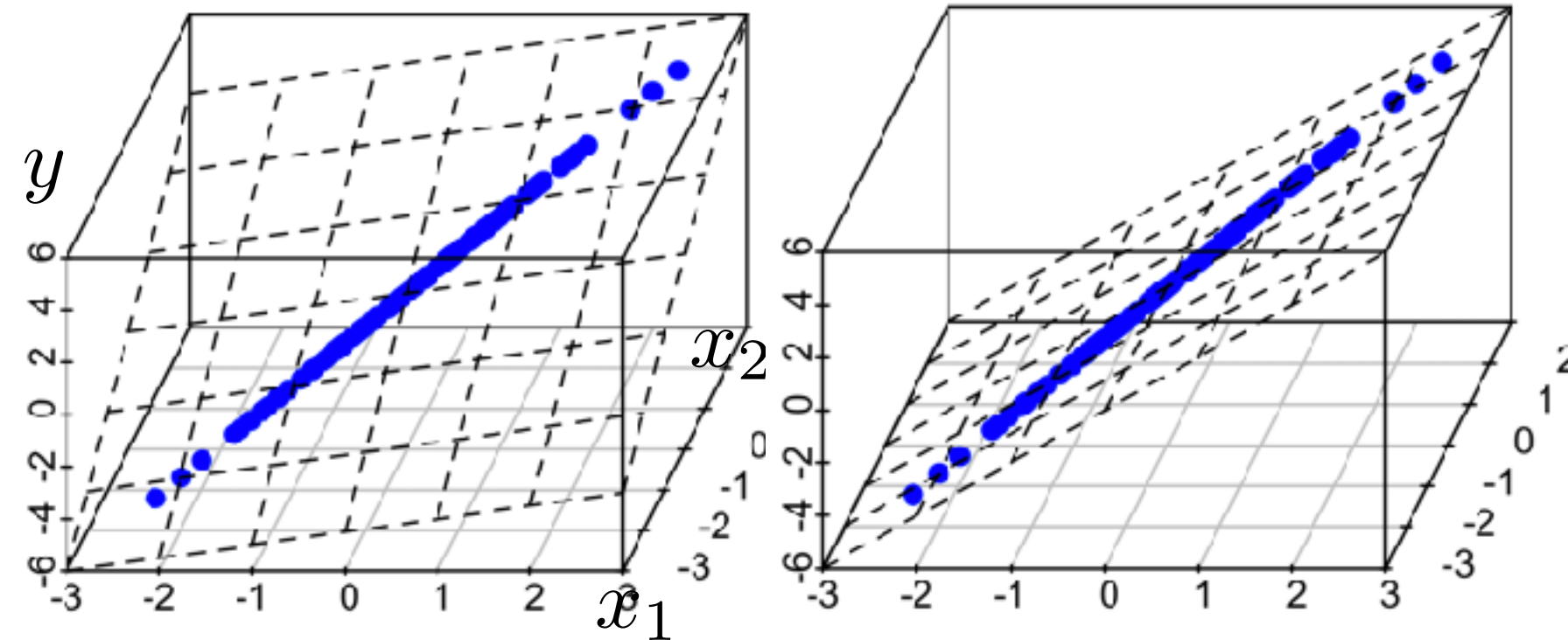
# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane



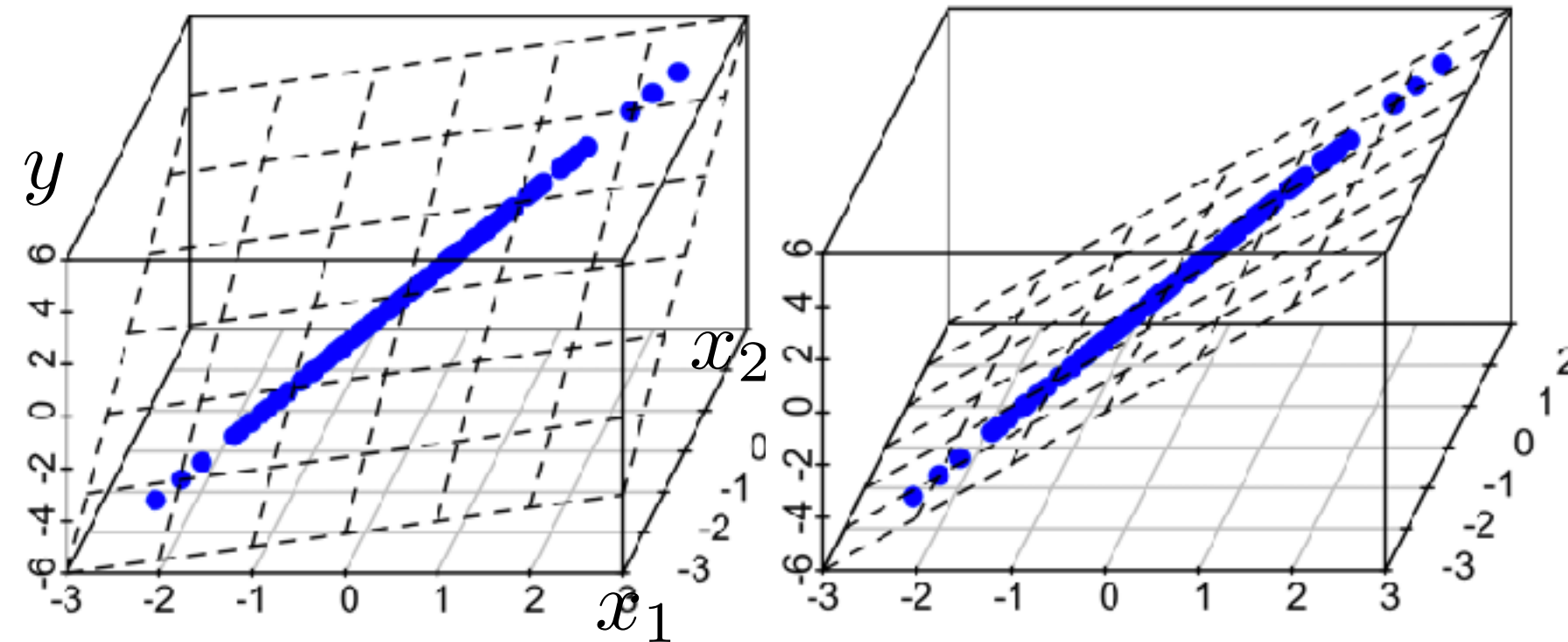
# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane
  - Then  $X^T X$  not invertible



# What can go wrong in practice?

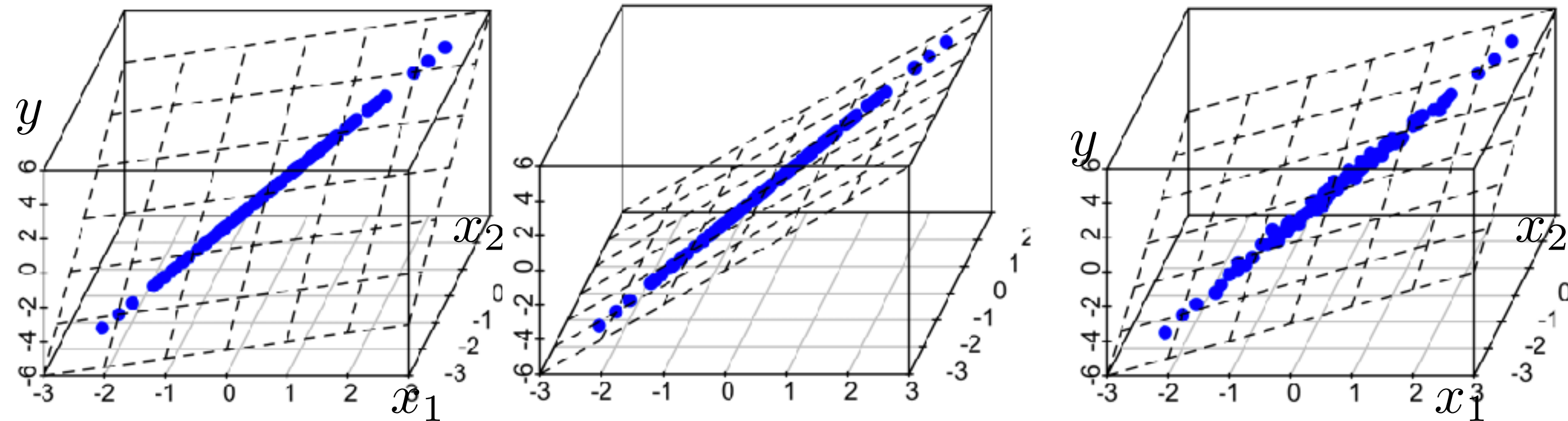
- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane
  - Then  $X^T X$  not invertible



- Sometimes there's technically a unique best hyperplane, but just because of noise

# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane
  - Then  $X^T X$  not invertible

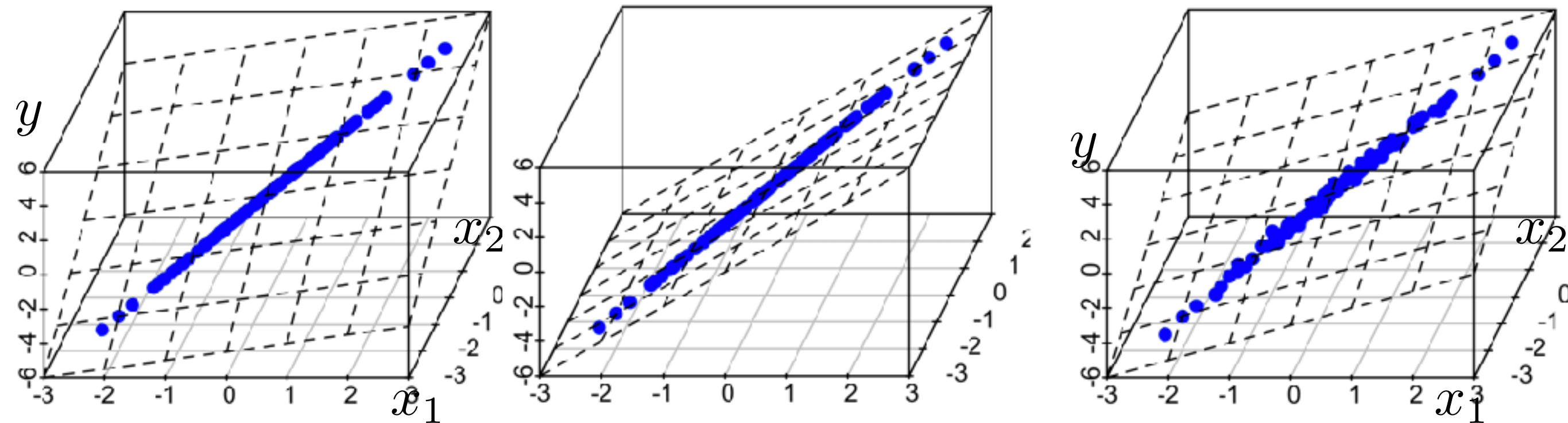


- Sometimes there's technically a unique best hyperplane, but just because of noise



# What can go wrong in practice?

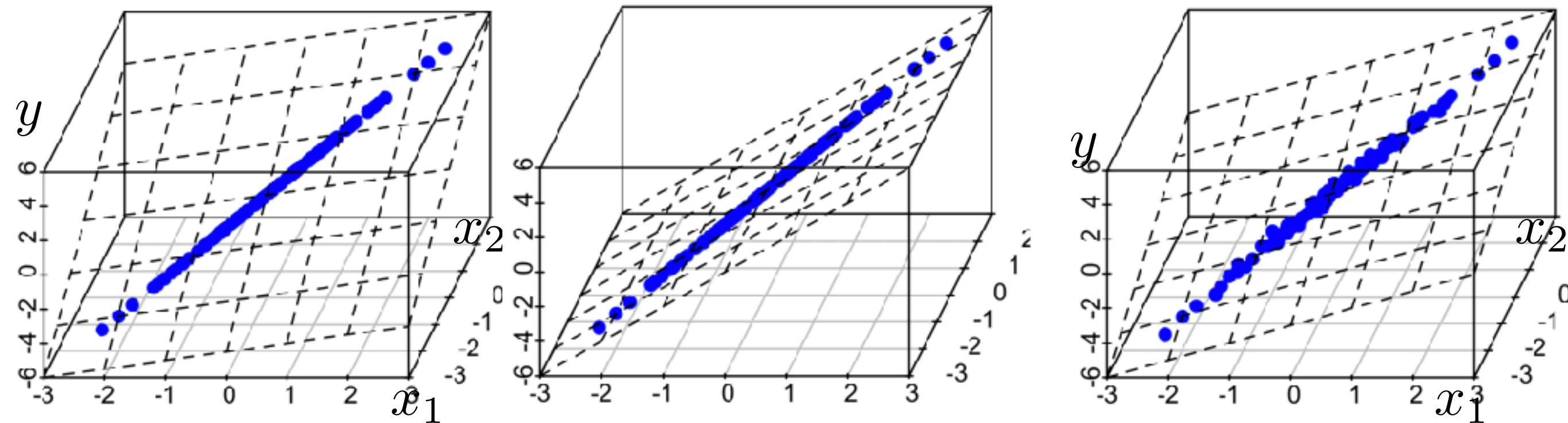
- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane
  - Then  $X^T X$  not invertible



- Sometimes there's technically a unique best hyperplane, but just because of noise
- Practical: real-life features often have this issue

# What can go wrong in practice?

- Does the linear regr. objective always “curve up”? No!
- Sometimes there isn't a unique best hyperplane
  - Then  $X^T X$  not invertible



- Sometimes there's technically a unique best hyperplane, but just because of noise
- Practical: real-life features often have this issue
- How to choose among hyperplanes? Preference for  $\theta$  components being near zero



# Regularizing linear regression

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2$$

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2$$

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2$$

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

What happens if  $\lambda < 0$  ?



# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

- When  $\lambda > 0$ , always “curves up” & can invert

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

- When  $\lambda > 0$ , always “curves up” & can invert
- Can also solve with an offset

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

- When  $\lambda > 0$ , always “curves up” & can invert
- Can also solve with an offset

- Can think of  $\lambda$  as hyperparameter of a learning algorithm

What happens if  $\lambda < 0$  ?

# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{\text{dxn, nxd}} + n \underbrace{\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

- When  $\lambda > 0$ , always “curves up” & can invert
- Can also solve with an offset

Exercise:  
write out  
the learning  
algorithm

What  
happens if  
 $\lambda < 0$  ?

- Can think of  $\lambda$  as hyperparameter of a learning algorithm



# Regularizing linear regression

- Linear regression with square penalty: ridge regression

$$J_{\text{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- Special case: ridge regression with no offset

$$J_{\text{ridge}}(\theta) = \frac{1}{n} (\tilde{X}\theta - \tilde{Y})^\top (\tilde{X}\theta - \tilde{Y}) + \lambda \|\theta\|^2$$

- Min at:  $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$

$$\Rightarrow \theta = (\underbrace{\tilde{X}^\top \tilde{X}}_{dxn, nxd} + n \underbrace{\lambda I}_{dxd})^{-1} \tilde{X}^\top \tilde{Y}$$

- When  $\lambda > 0$ , always “curves up” & can invert
- Can also solve with an offset

Exercise:  
write out  
the learning  
algorithm

What  
happens if  
 $\lambda < 0$  ?

- Can think of  $\lambda$  as hyperparameter of a learning algorithm
- How to choose  $\lambda$ ? One option: cross validation (see HW!)