

<https://introml.mit.edu/>

# 6.390 Intro to Machine Learning

## Lecture 2: Regularization and Cross-validation

Shen Shen

Feb 9, 2026

3pm, Room 10-250

[Slides and Lecture Recording](#)

Recall

Let

$$X = \begin{bmatrix} \text{---} \mathbf{x}^{(1)\top} \text{---} \\ \vdots \\ \text{---} \mathbf{x}^{(n)\top} \text{---} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^{(1)} & \cdots & \mathbf{x}_d^{(1)} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^{(n)} & \cdots & \mathbf{x}_d^{(n)} \end{bmatrix} \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$\in \mathbb{R}^{n \times d}$                        $\in \mathbb{R}^{n \times 1}$                        $\in \mathbb{R}^{d \times 1}$

Then

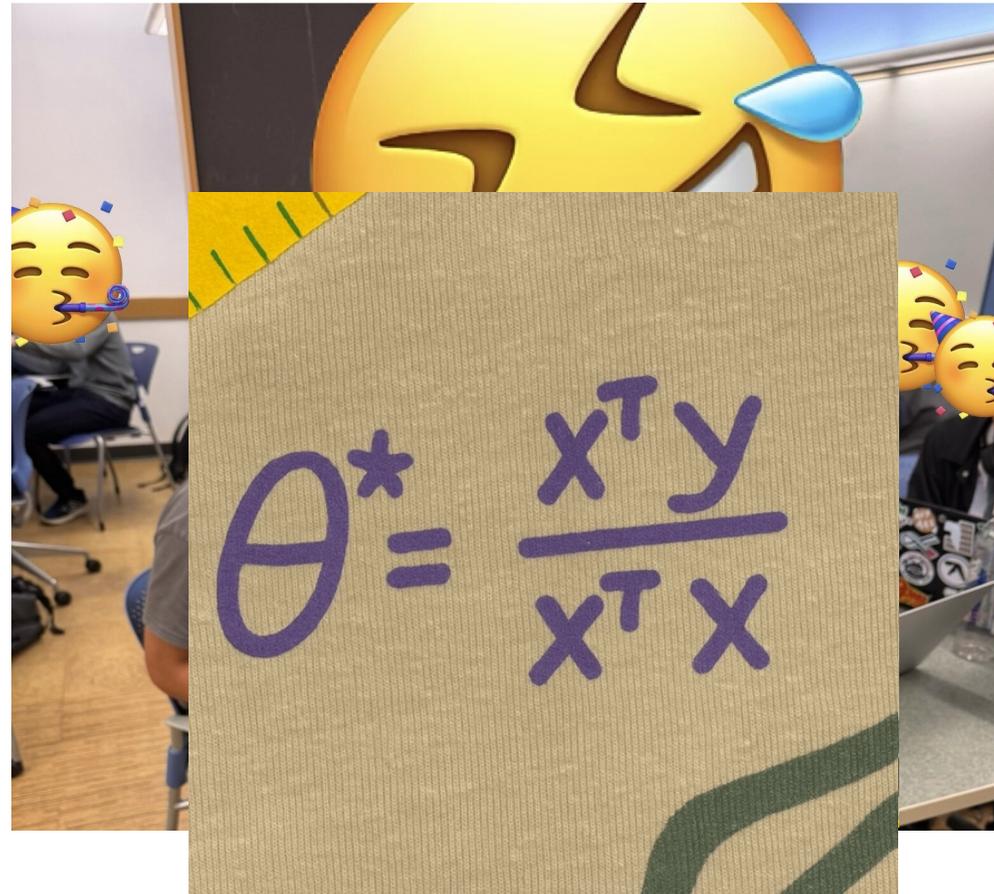
$$J(\theta) = \frac{1}{n} (X\theta - Y)^\top (X\theta - Y) \in \mathbb{R}^{1 \times 1}$$

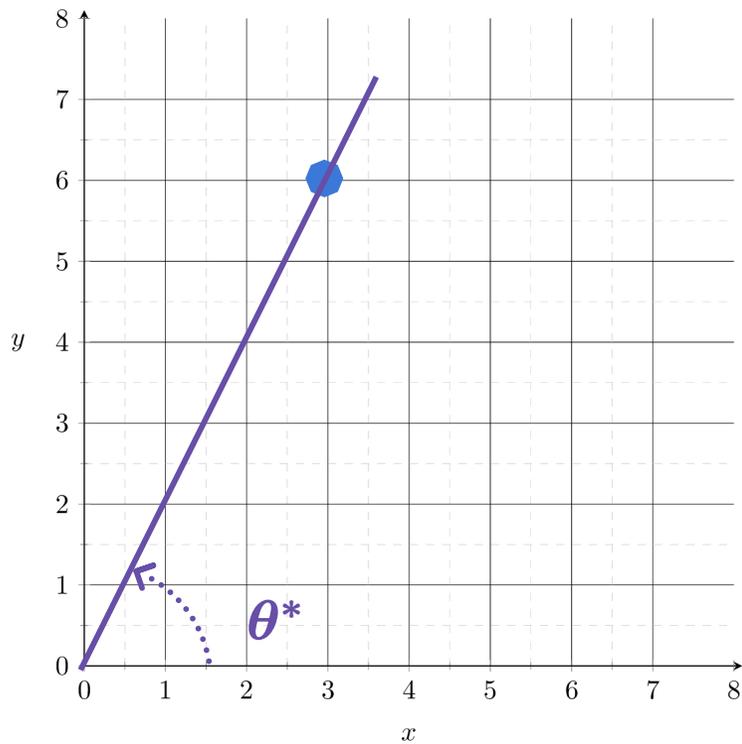
By matrix calculus and optimization

$$\theta^* = (X^\top X)^{-1} X^\top Y \in \mathbb{R}^{d \times 1}$$

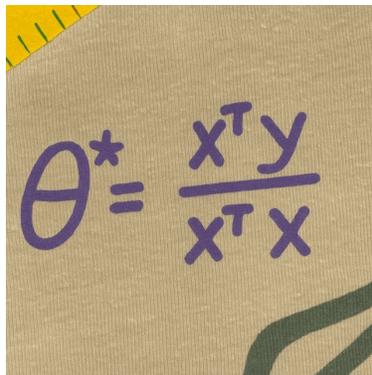
$$\theta^* = (X^T X)^{-1} X^T Y$$

Jane street shirt





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



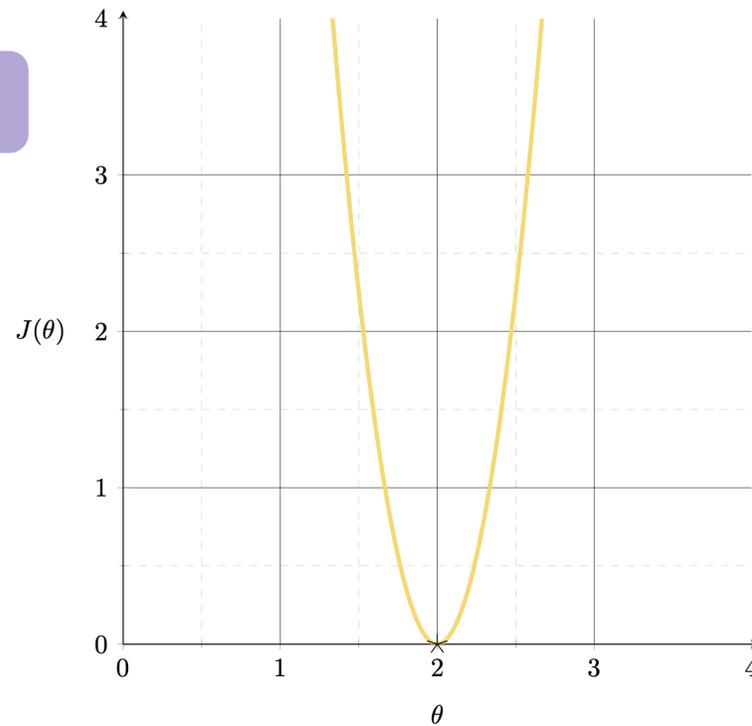
$$X = [3]$$

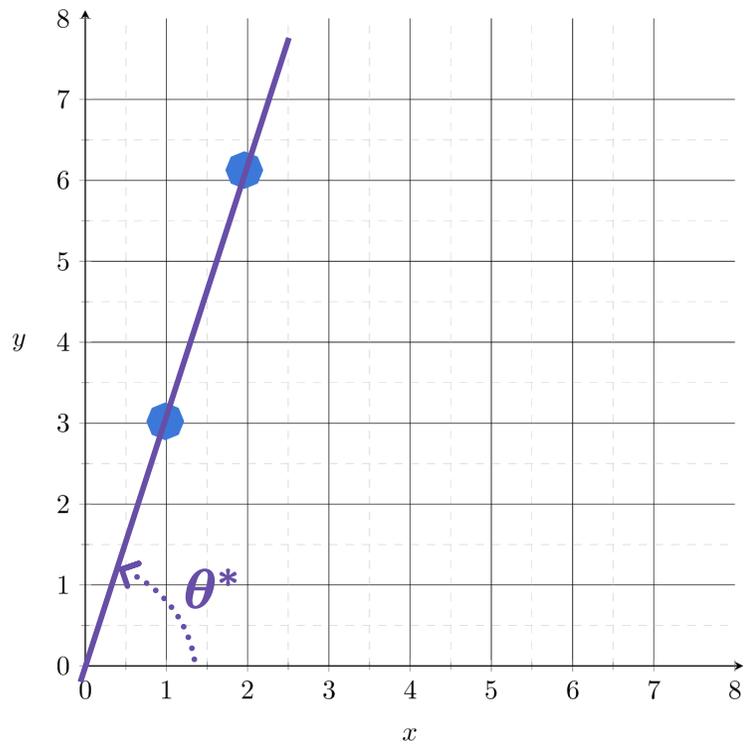
$$Y = [6]$$

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

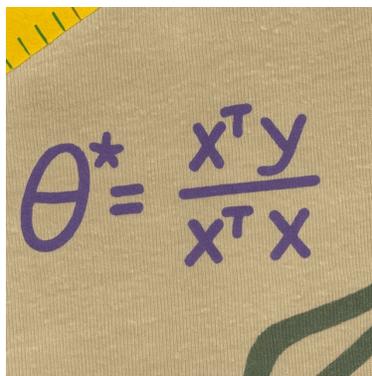
$$= \frac{X^T Y}{X^T X} = \frac{6}{3} = 2$$

$$J(\theta) = (3\theta - 6)^2$$



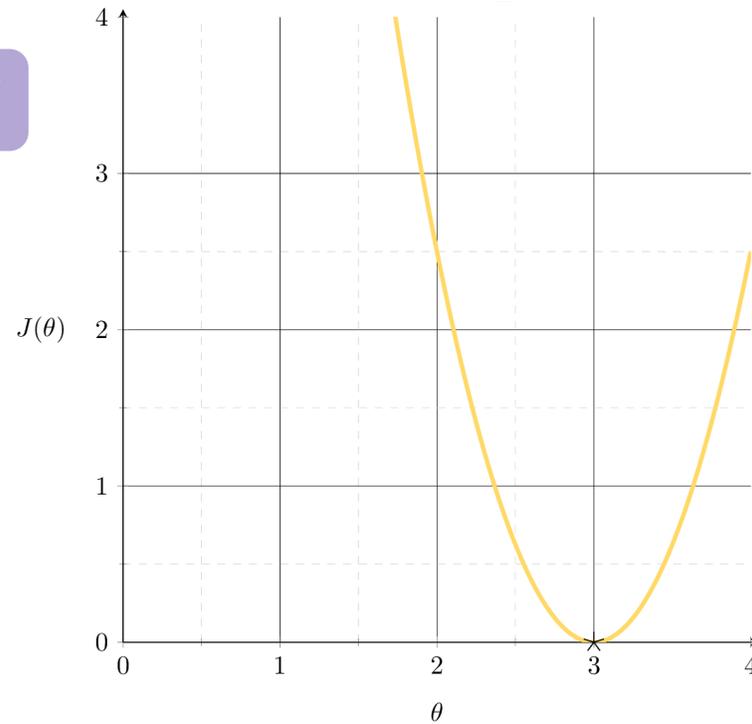


$$\theta^* = (X^T X)^{-1} X^T Y$$



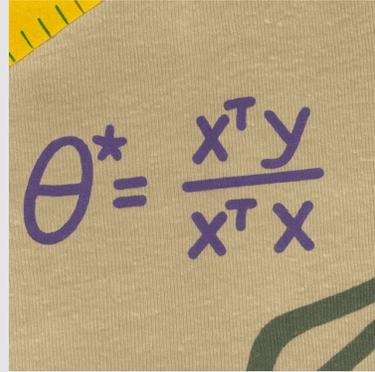
$$X = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad Y = \begin{bmatrix} 3 \\ 6 \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} [(\theta - 3)^2 + (2\theta - 6)^2]$$



$$\theta^* = (X^T X)^{-1} (X^T Y) = \frac{X^T Y}{X^T X} = \frac{15}{5} = 3$$

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



A photograph of a piece of textured, light-brown fabric with a yellow and green patterned border. The equation  $\theta^* = \frac{X^T y}{X^T X}$  is written in purple ink on the fabric.

<https://shenshen.mit.edu/demos/ridge/d2-unique-solution.html>

$$X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad Y = \begin{bmatrix} 2 \\ 3 \\ 5 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad X^T Y = \begin{bmatrix} 7 \\ 8 \end{bmatrix}$$

$$\theta^* = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 7 \\ 8 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

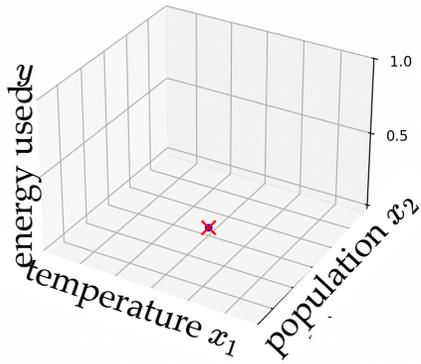
# Outline

- The "*trouble*" with the closed-form solution
  - visually, practically, mathematically
- Regularization and ridge regression
- Cross-validation

<https://shenshen.mit.edu/demos/ridge/n-less-d-interactive.html?embed>

<https://shenshen.mit.edu/demos/ridge/collinear-interactive.html?embed>

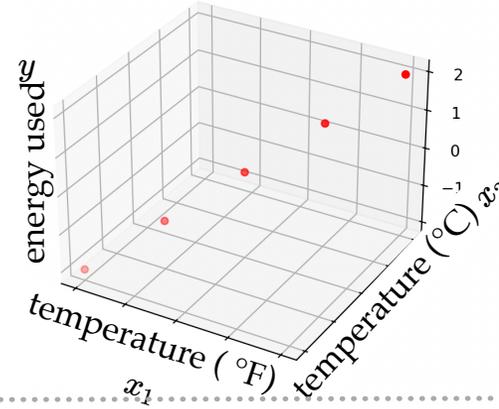
data



(a)  $n < d$

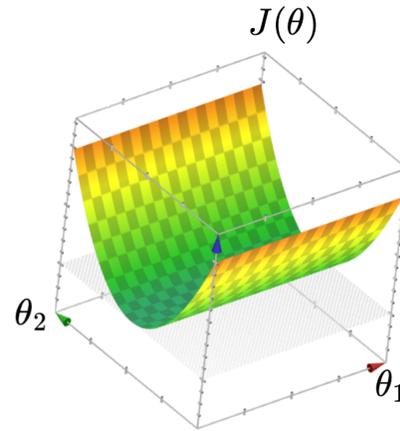
e.g. genomics, NLP

(b) Linearly-dependent features:



e.g. temp  $^{\circ}$ F/ $^{\circ}$ C,  
age/birth\_year, ...

MSE



closed-form formula

$$\theta^* = (X^T X)^{-1} X^T Y \quad \text{undefined}$$

$\infty$  many optimal  $\theta^*$

optimal solution

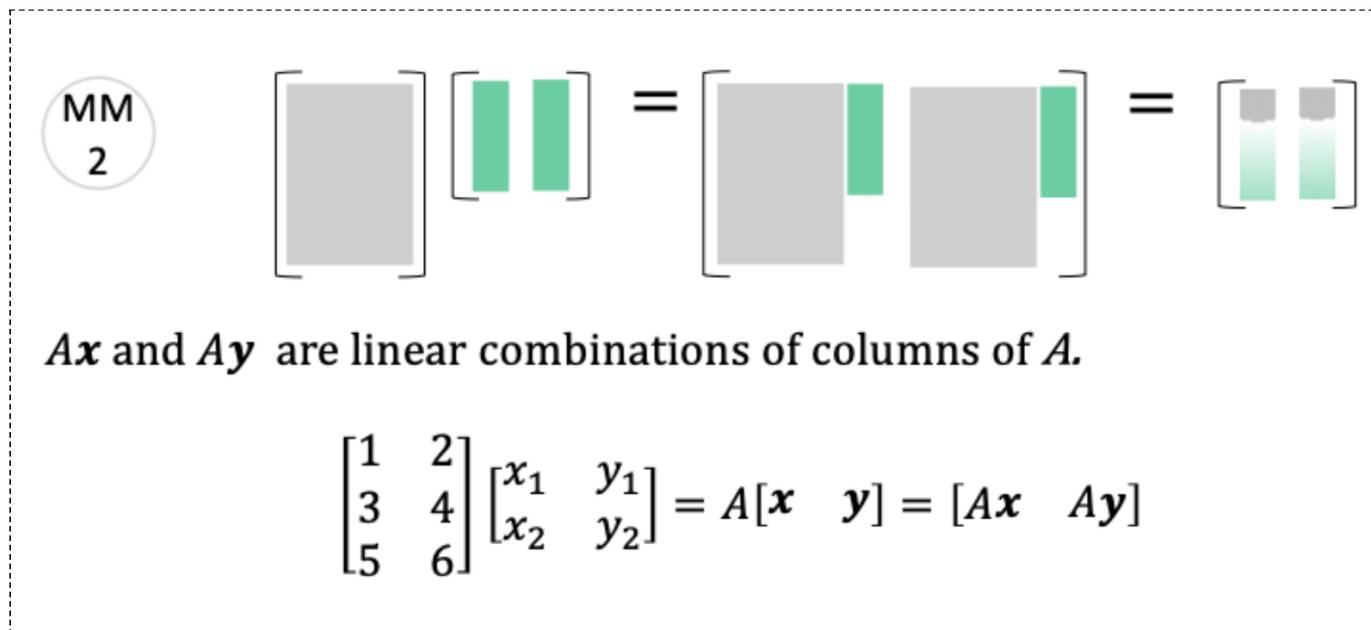
Not enough info to pin down a unique solution

mathematically,

$(X^\top X)$  is singular



$X$  is not full column rank

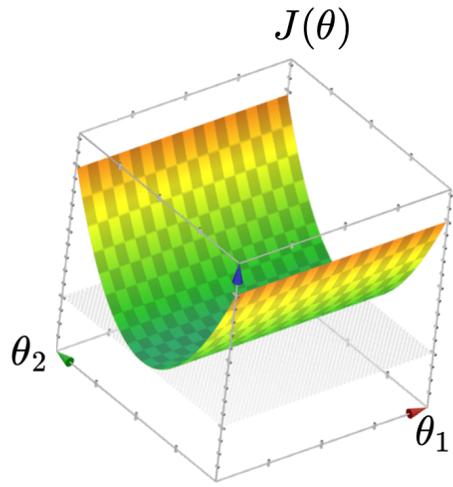


demo (a)  $n < d$ : 1 sample, 2 features

$$X^\top X = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \begin{bmatrix} 2 & 3 \end{bmatrix} = \begin{bmatrix} 4 & 6 \\ 6 & 9 \end{bmatrix}$$

demo (b) Collinear:  $x_2 = 1.5 \cdot x_1$

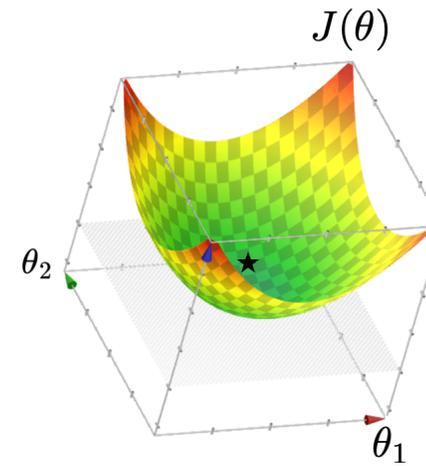
$$X^\top X = \begin{bmatrix} 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 4 & 6 \\ 6 & 9 \end{bmatrix} = \begin{bmatrix} 56 & 84 \\ 84 & 126 \end{bmatrix}$$



When  $X$  is not full column rank

- $J(\theta)$  has a "flat" bottom
- This 🙅 formula is not well-defined
- Infinitely many optimal hyperplanes

formula isn't wrong, data is trouble-making 🙄



Typically,  $X$  is full column rank

- $J(\theta)$  "curves up" everywhere
- $\theta^* = (X^\top X)^{-1} X^\top Y$
- unique optimal hyperplane

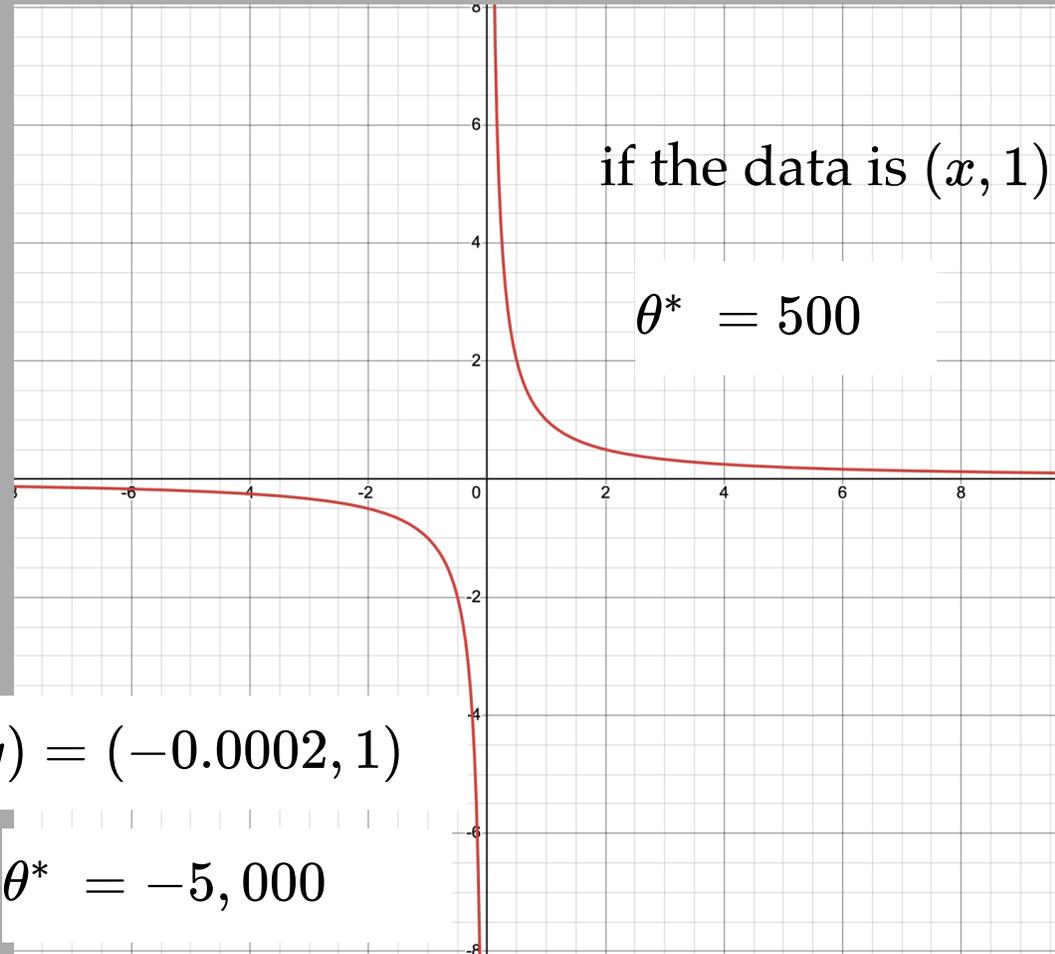
$X^\top X$  is "more" invertible 😊

minimum eigenvalue of  $(X^\top X)$  increasing

$$\theta^* = \frac{x^T y}{x^T x}$$

assume  $n = 1$  and  $y = 1$

then  $\theta^* = \frac{1}{x}$



<https://shenshen.mit.edu/demos/ridge/sensitive-ols.html>

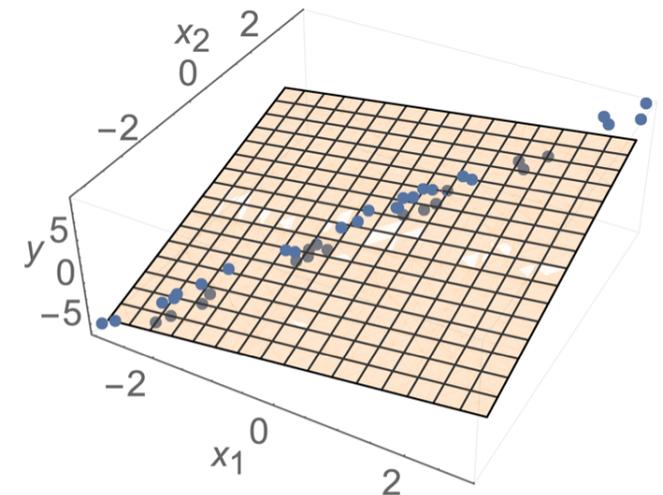
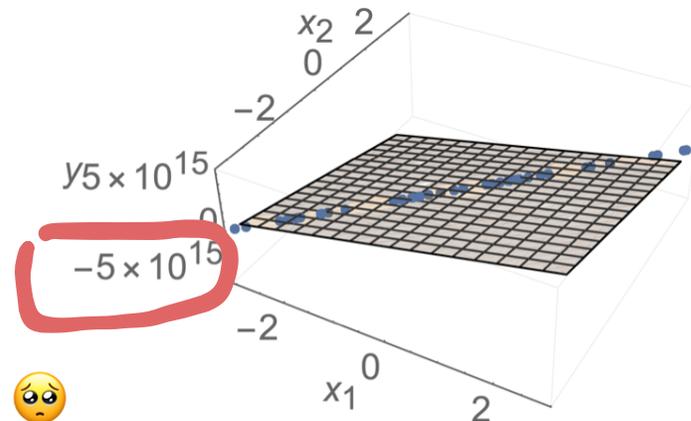
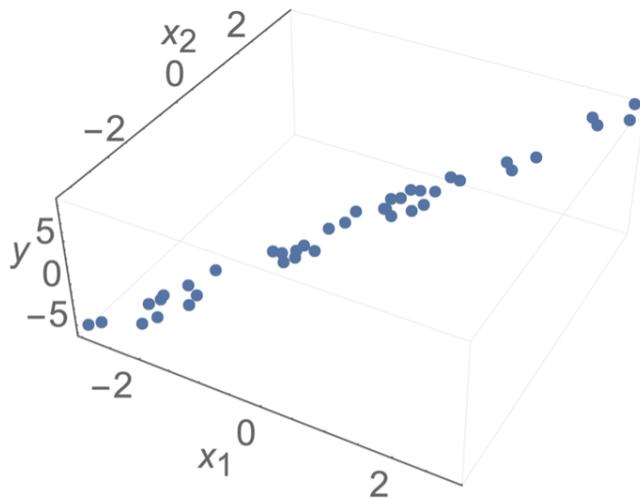
when  $X^T X$  is *almost* singular

technically,  $\theta^* = (X^T X)^{-1} X^T Y$  exists and gives the unique optimal hyperplane

practically,  $\theta^*$  tends to have huge magnitude

$\theta^*$  tends to be very sensitive to the small changes in the data

lots of other  $\theta$ s fit the training data almost equally well



# Outline

- The "*trouble*" with the closed-form solution
- Regularization and ridge regression
  - hyperparameters
- Cross-validation

# Ridge Regression: Objective

- Many  $\theta$  give similar loss, but some have huge magnitude... unstable!
  - small change in  $x \Rightarrow$  wildly different prediction
- Idea: penalize large  $\theta$  in our objective, a.k.a. (explicit) *regularization*
- Ridge objective:

$$J_{\text{ridge}}(\theta) = \underbrace{\frac{1}{n} (X\theta - Y)^\top (X\theta - Y)}_{\text{MSE (on training data)}} + \underbrace{\lambda \|\theta\|^2}_{\text{penalty (on parameter magnitude)}}$$

$\lambda > 0$  controls how heavily we penalize magnitude relative to MSE

<https://shenshen.mit.edu/demos/ridge/ridge-lambda.html?embed>

# Ridge Regression: Solution

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

$$\theta_{\text{ridge}}^* = (X^\top X + n\lambda I)^{-1} X^\top Y$$

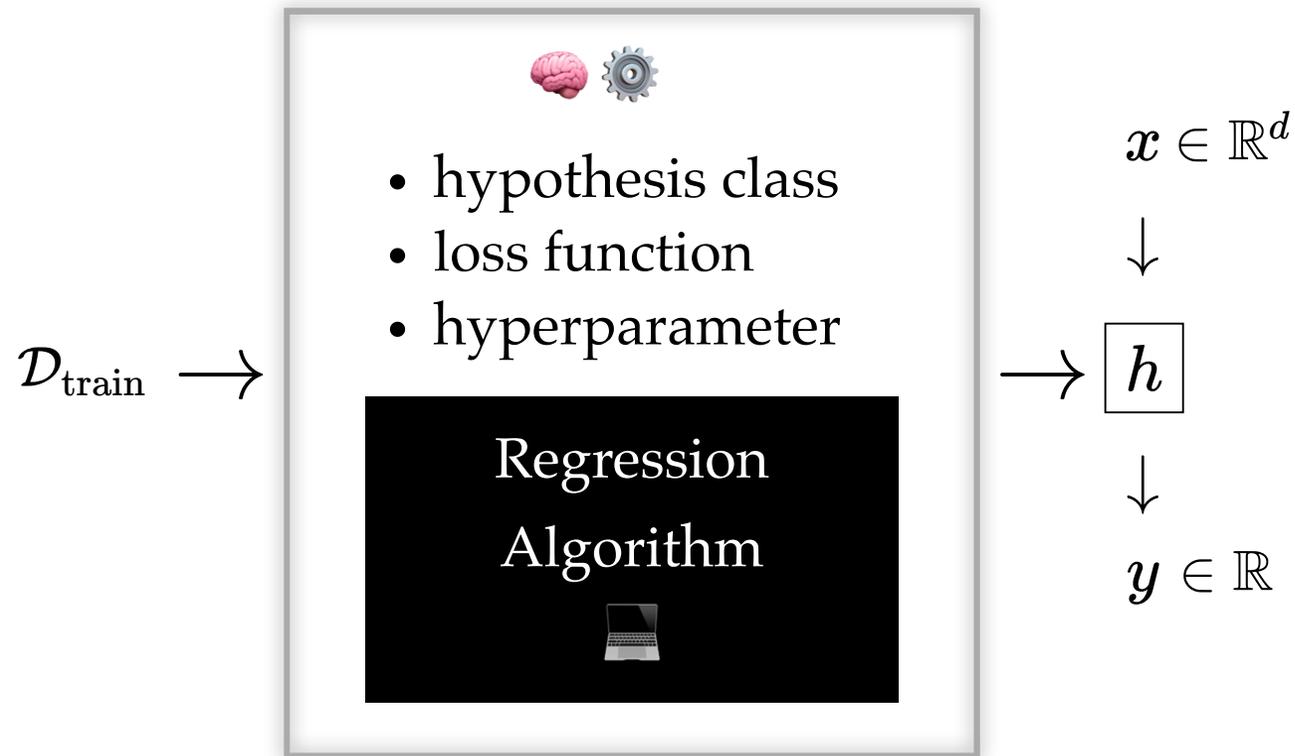
for  $\lambda > 0$ :  $X^\top X + n\lambda I$  is always invertible, so  $\theta_{\text{ridge}}^*$  always exists and is unique

How does  $\lambda$  affect the learned  $\theta$ ?

- $\lambda = 0$ ? No penalty — reduces to OLS
- $\lambda = 1000$ ? Huge penalty — forces  $\theta \approx 0$
- $\lambda = -100$ ? *Rewards* large  $\theta$  — counter-productive!

this is why we require  $\lambda > 0$

$\lambda$  is a hyperparameter



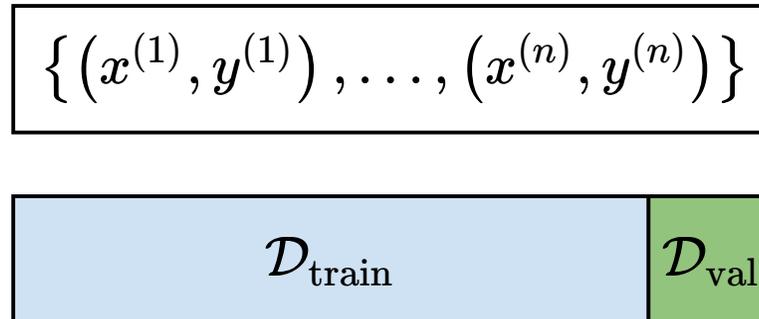
- affects learning outcome, and not learned by the algorithm
- we already saw a hyperparameter (in lab 1, how many random regressor tried)

# Outline

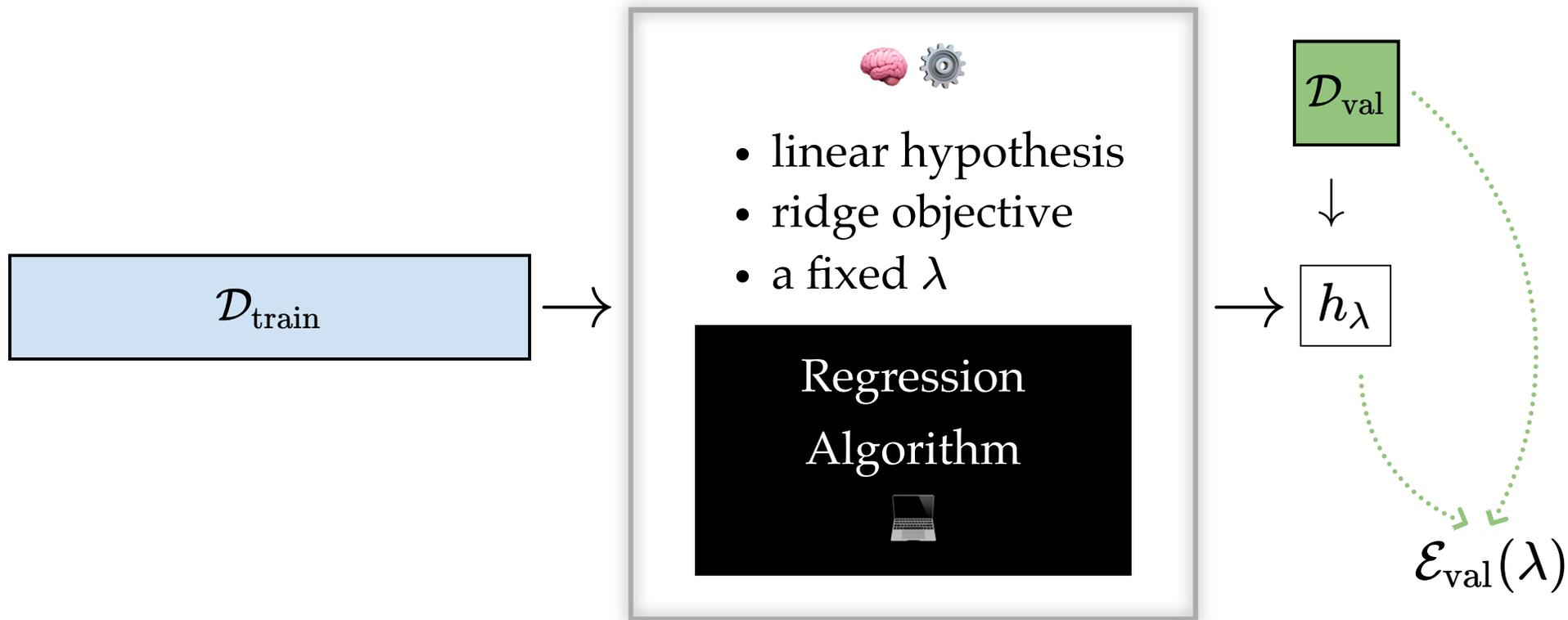
- The "*trouble*" with the closed-form solution
- Regularization and ridge regression
- Cross-validation

We need to choose hyperparameters (like  $\lambda$ )

- Can't use training error (would always pick  $\lambda = 0$ )
- Can't use test error (we don't have test data)
  
- Hold-out some data



- Use  $\mathcal{D}_{\text{val}}$  to evaluate how good a hyperparameter  $\lambda$  is



**for each**  $\lambda \in \{0.1, 1, 10\}$ :

train on  $\mathcal{D}_{\text{train}}$  with  $\lambda$

compute  $\mathcal{E}_{\text{val}}(\lambda)$  on  $\mathcal{D}_{\text{val}}$

**return**  $\arg \min_{\lambda} \mathcal{E}_{\text{val}}(\lambda)$

in this example, compare  $\mathcal{E}_{\text{val}}(.1)$ ,  $\mathcal{E}_{\text{val}}(1)$ , and  $\mathcal{E}_{\text{val}}(10)$ ,  
return the  $\lambda$  corresponding to smallest validation error

# Cross-validation

**for each**  $\lambda \in \{0.1, 1, 10\}$ :

**for**  $i = 1, \dots, 5$ :

train  $h_i$  on  $\mathcal{D} \setminus \mathcal{D}_i$  with  $\lambda$

$\mathcal{E}_i$  = error on  $\mathcal{D}_i$

$$\mathcal{E}_{\text{val}}(\lambda) = (\mathcal{E}_1 + \dots + \mathcal{E}_5)/5$$

**return**  $\lambda^* = \arg \min_{\lambda} \mathcal{E}_{\text{val}}(\lambda)$

outer loop of  $\lambda \in \{0.1, 1, 10\}$ :



How many hypotheses trained in this example to pick  $\lambda^*$ ?

$$\mathcal{E}_{\text{val}}(\lambda) = (\mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3 + \mathcal{E}_4 + \mathcal{E}_5)/5$$

$$\theta_{\text{final}}^* = (X^T X + n\lambda^* I)^{-1} X^T Y$$

finally train using the chosen  $\lambda^*$ ,  
and data from all of  $\mathcal{D}$

# Summary

- When  $X^\top X$  is singular or ill-conditioned, OLS is undefined or overfits.
- Regularization combats overfitting by penalizing large  $\theta$ .
- Ridge regression adds  $\lambda\|\theta\|^2$  to the objective — still has a closed-form solution.
- $\lambda$  is a hyperparameter that trades off fit vs. regularization.
- Validation and cross-validation provide principled ways to choose  $\lambda$ .

# Cross-validation

**for each**  $\lambda \in \{0.1, 1, 10\}$ :

**for**  $i = 1, \dots, 5$ :

train  $h_i$  on  $\mathcal{D} \setminus \mathcal{D}_i$  with  $\lambda$

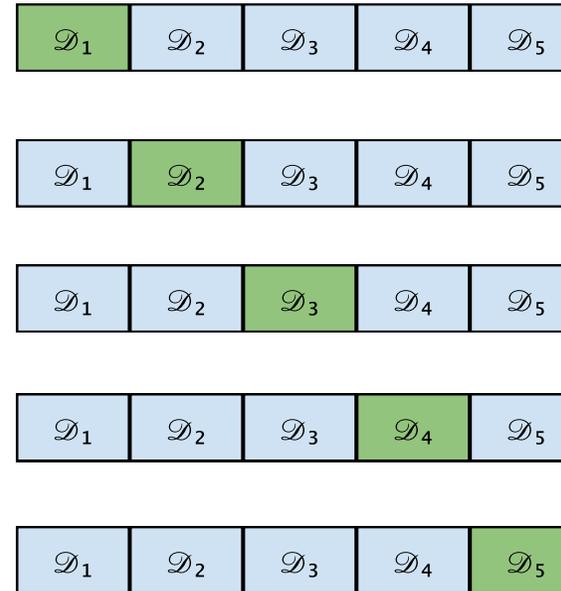
$\mathcal{E}_i =$  error on  $\mathcal{D}_i$

$$\mathcal{E}_{\text{val}}(\lambda) = (\mathcal{E}_1 + \dots + \mathcal{E}_5)/5$$

**return**  $\lambda^* = \arg \min_{\lambda} \mathcal{E}_{\text{val}}(\lambda)$

$$\theta_{\text{final}}^* = (X^{\top} X + n\lambda^* I)^{-1} X^{\top} Y$$

outer loop of  $\lambda \in \{0.1, 1, 10\}$ :



$$\mathcal{E}_{\text{val}}(\lambda) = (\mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3 + \mathcal{E}_4 + \mathcal{E}_5)/5$$

How many hypotheses trained in this example to pick  $\lambda^*$ ?

finally train using data from all of  $\mathcal{D}$