

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

6.390 Intro to Machine Learning

Lecture 2: Regularization and Cross-validation

Shen Shen

Sept 11, 2025

11am, Room 10-250

[Interactive Slides and Lecture Recording](#)

Outline

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

Recall

- Linear hypothesis class:

$$h(x; \theta) = \theta^T x$$

A diagram illustrating the components of the linear hypothesis equation $h(x; \theta) = \theta^T x$. The word "parameters" is positioned below θ , and the word "features" is positioned below x . Dotted curved arrows point from each word to its corresponding variable in the equation.

- Squared loss function:

$$\mathcal{L}(h(x^{(i)}), y^{(i)}) = (\theta^T x^{(i)} - y^{(i)})^2$$

A diagram illustrating the components of the squared loss function equation $\mathcal{L}(h(x^{(i)}), y^{(i)}) = (\theta^T x^{(i)} - y^{(i)})^2$. The word "loss" is positioned below \mathcal{L} . The words "guess (prediction)" and "label" are positioned below $\theta^T x^{(i)}$ and $y^{(i)}$ respectively. Dotted curved arrows point from "loss" to \mathcal{L} , from "guess (prediction)" to $\theta^T x^{(i)}$, and from "label" to $y^{(i)}$.

See lec1/rec1 for discussion of the offset.

Recall

Let

$$X = \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times d} \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times 1} \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \in \mathbb{R}^{d \times 1}$$

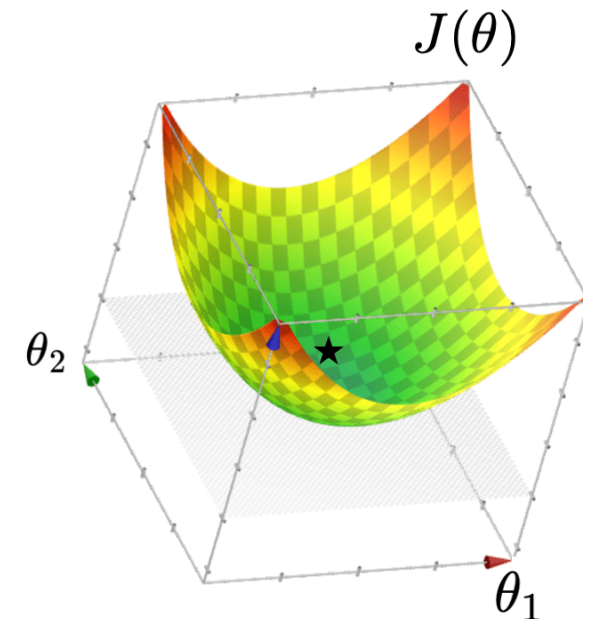
Then

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

By matrix calculus and optimization

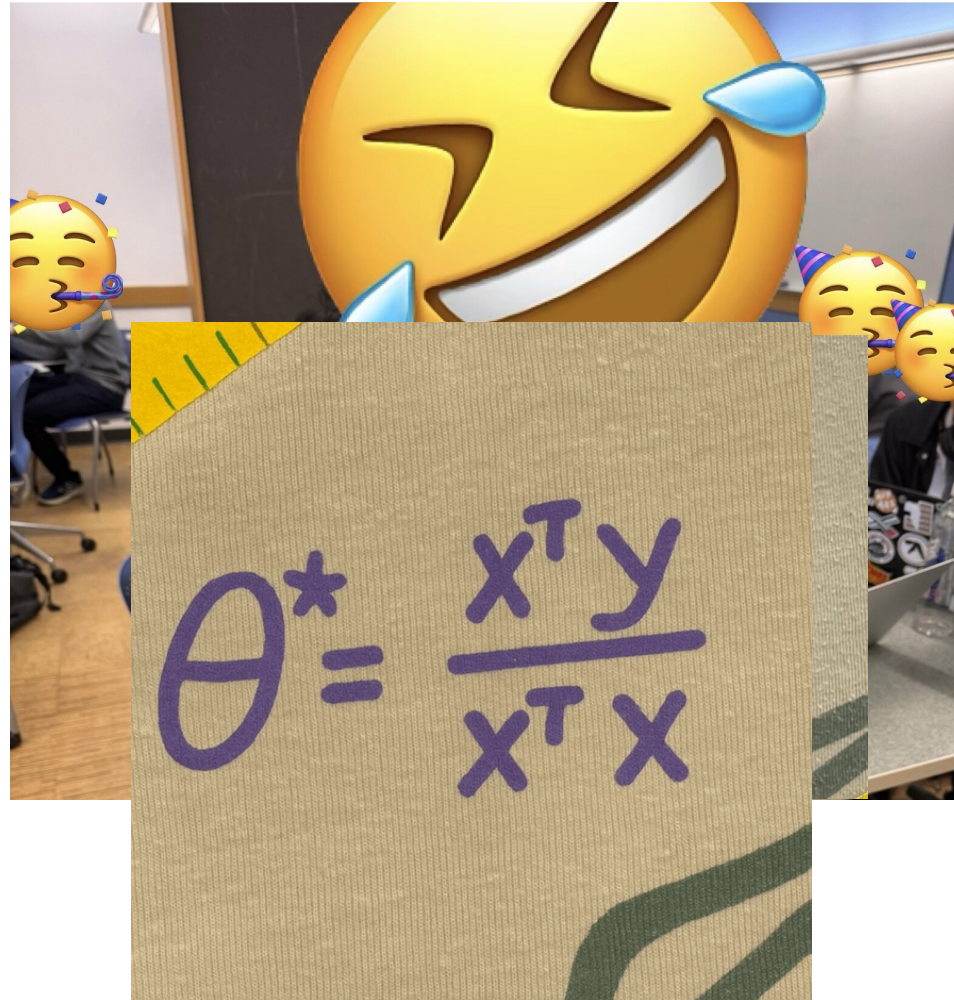
$$\theta^* = (X^\top X)^{-1} X^\top Y$$

$$X^\top X \in \mathbb{R}^{d \times d} \quad X^\top Y \in \mathbb{R}^{d \times 1}$$

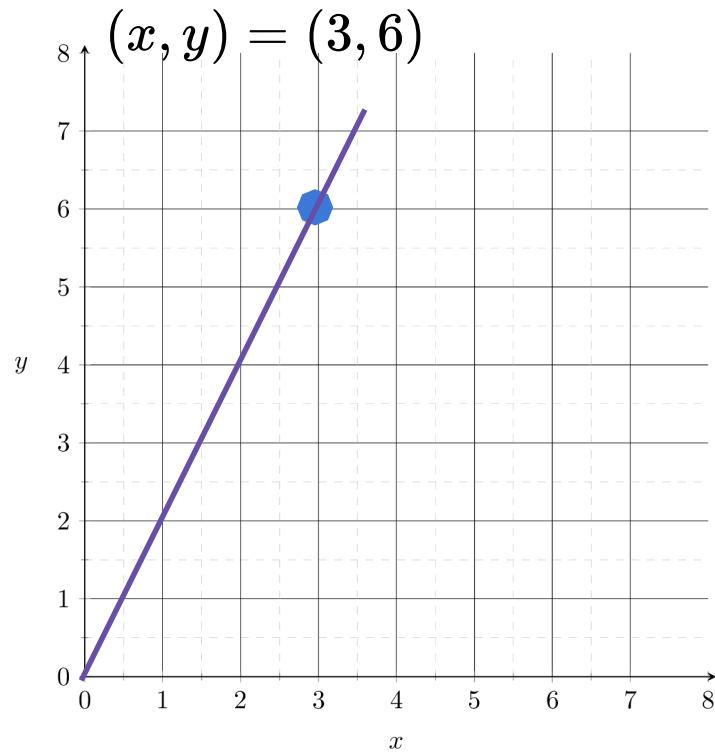


Spotted in lab:

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



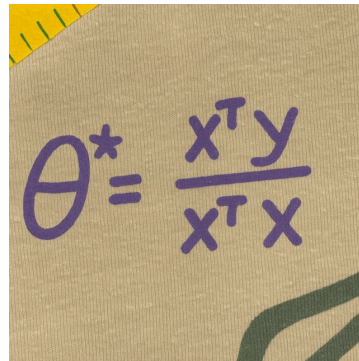
1d-feature training data



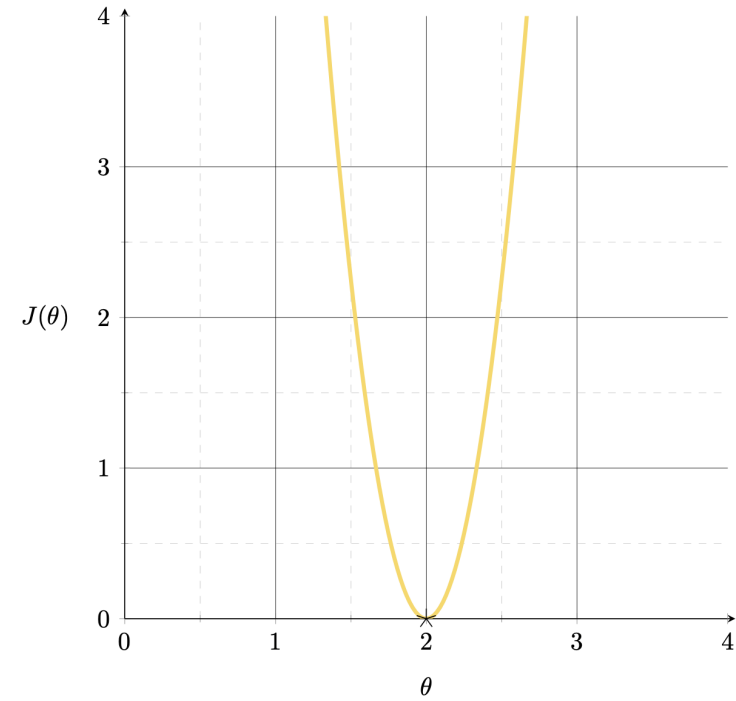
$$\theta^* = (X^\top X)^{-1} X^\top Y$$

$$X = x = [3]$$

$$Y = y = [6]$$


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

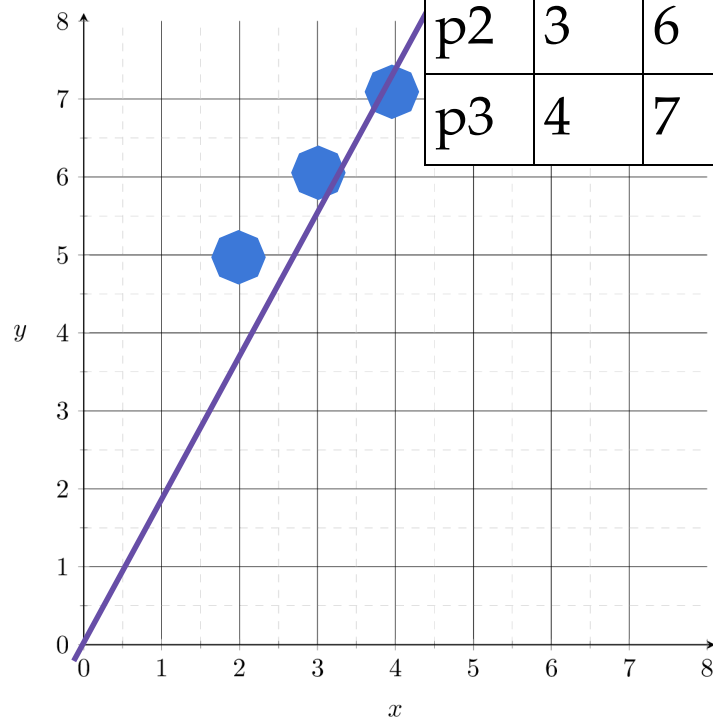
$$J(\theta) = (3\theta - 6)^2$$



$$\theta^* = (xx)^{-1}(xy) = \frac{xy}{xx} = \frac{y}{x} = \frac{6}{3} = 2$$

1-d feature
training data set

	x	y
p1	2	5
p2	3	6
p3	4	7



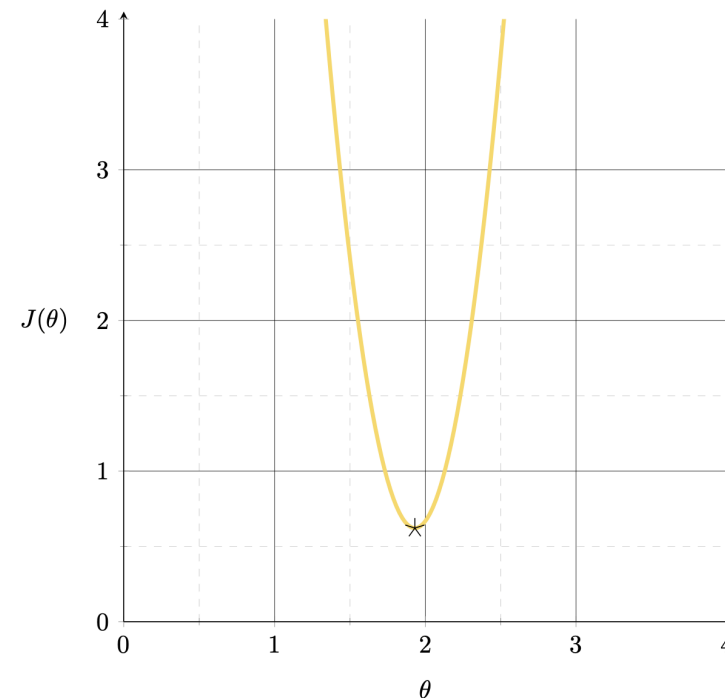
$$\theta^* = (X^T X)^{-1} X^T Y$$

$$X = \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$$

$$Y = \begin{bmatrix} 5 \\ 6 \\ 7 \end{bmatrix}$$

A photograph of a piece of brown paper with the formula $\theta^* = \frac{X^T Y}{X^T X}$ handwritten in purple ink.

$$J(\theta) = \frac{1}{3} [(2\theta - 5)^2 + (3\theta - 6)^2 + (4\theta - 7)^2]$$



$$\theta^* = \left(\begin{bmatrix} 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} \right)^{-1} \begin{bmatrix} 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \\ 7 \end{bmatrix} = \frac{X^T Y}{X^T X} = \frac{56}{29} \approx 1.93$$

Outline

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

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

$$d = 1$$

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

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

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

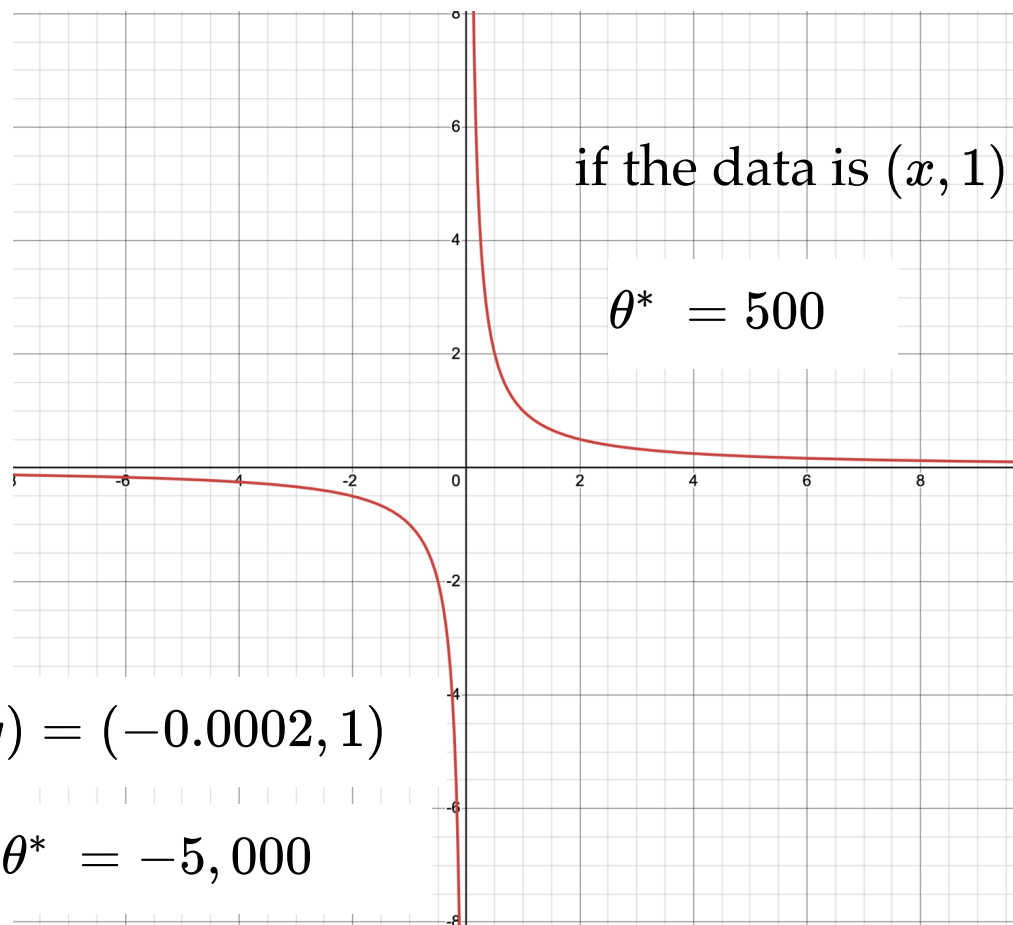
most of the time, behaves nicely

if the data is $(x, 1) = (0.002, 1)$

$$\theta^* = 500$$

if the data is $(x, y) = (-0.0002, 1)$

$$\theta^* = -5,000$$



more generally, $d \geq 1$

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

most of the time, behaves nicely

but run into trouble when $(X^\top X)$ is singular



$(X^\top X)$ has zero
eigenvalue(s)



the determinant of
 $(X^\top X)$ is zero



$(X^\top X)$ is not full rank



X is not full column rank

if X is not full column rank, then $X^\top X$ is singular

MM
2

$$\begin{bmatrix} \text{gray} \end{bmatrix} \begin{bmatrix} \text{green} & \text{green} \end{bmatrix} = \begin{bmatrix} \text{gray} & \text{green} & \text{gray} & \text{green} \end{bmatrix} = \begin{bmatrix} \text{green} & \text{green} \end{bmatrix}$$

$A\mathbf{x}$ and $A\mathbf{y}$ are linear combinations of columns of A .

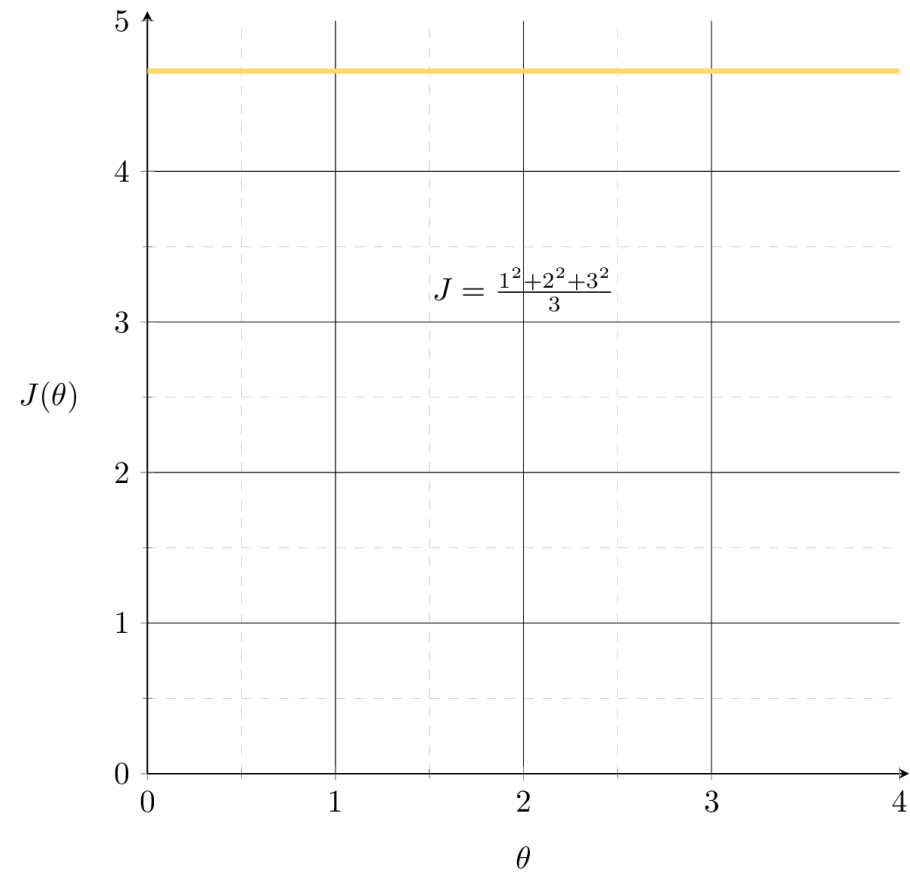
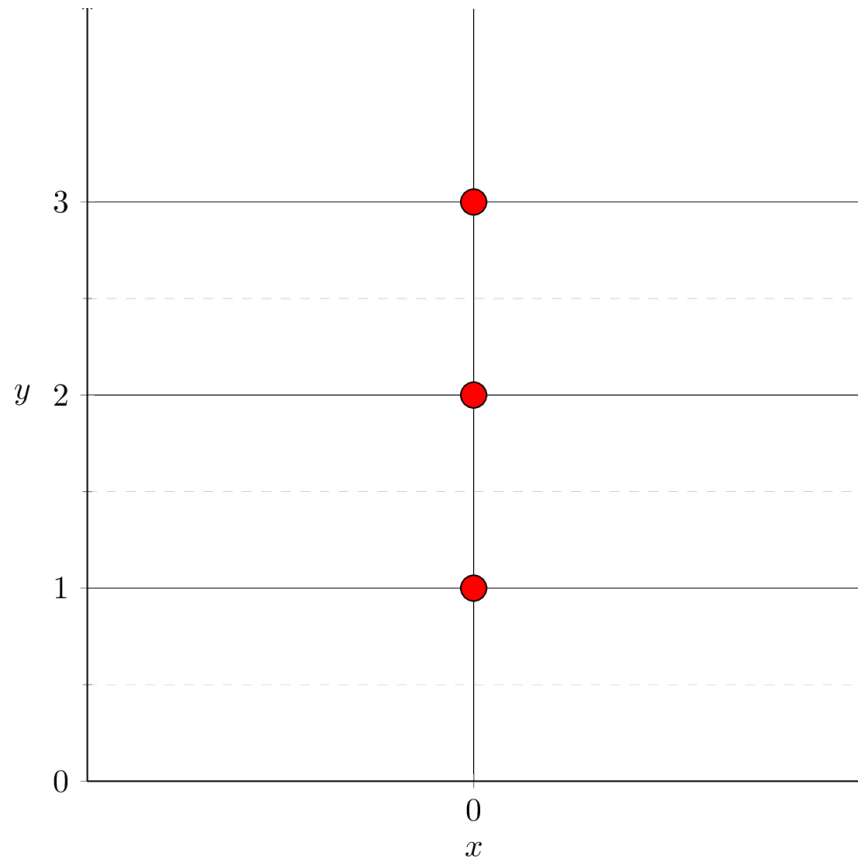
$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = A \begin{bmatrix} \mathbf{x} & \mathbf{y} \end{bmatrix} = \begin{bmatrix} A\mathbf{x} & A\mathbf{y} \end{bmatrix}$$

X is not full column rank when:

- a. $d = 1$ and $X \in \mathbb{R}^{n \times 1}$ is simply an all-zero vector, or
- b. $n < d$, or
- c. columns (features) in X are linearly dependent.

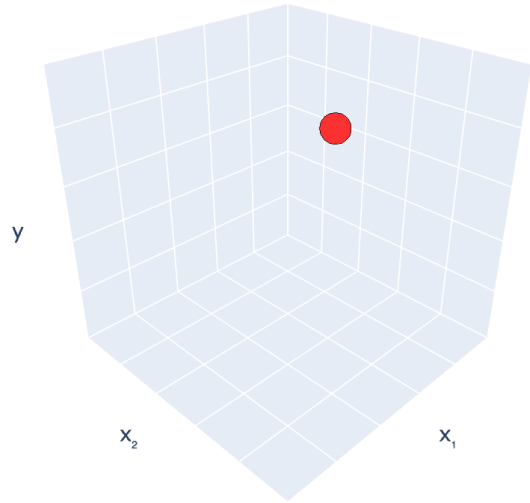
all three cases have similar visual interpretations

(a). $d = 1$ and $X \in \mathbb{R}^{n \times 1}$ is simply an all-zero vector



infinitely many optimal θ

(b). $n < d$

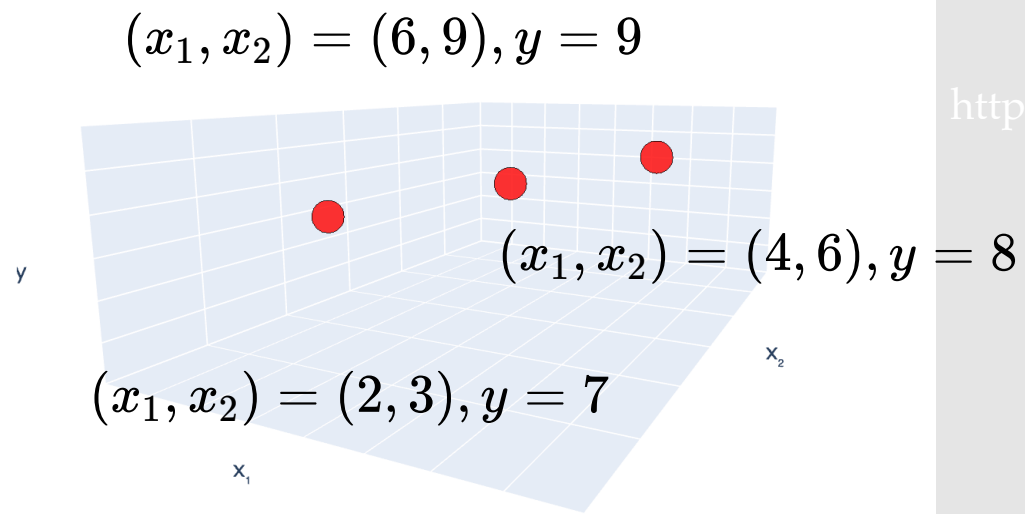


$$(x_1, x_2) = (2, 3), y = 4$$

<https://shenshen.mit.edu/demos/ridge/n>

infinitely many optimal θ

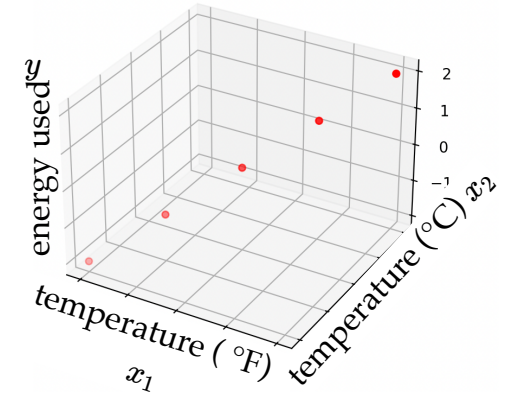
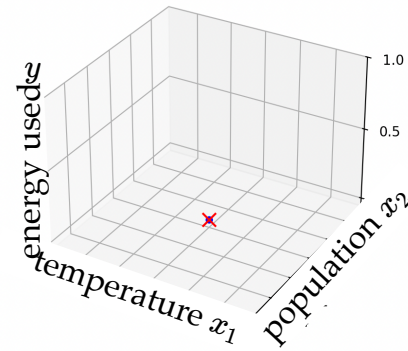
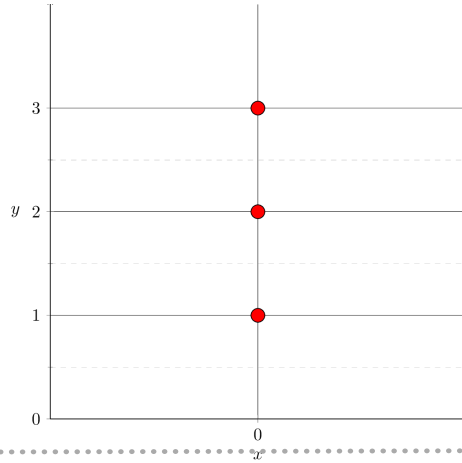
(c). columns (features) in X are linearly dependent.



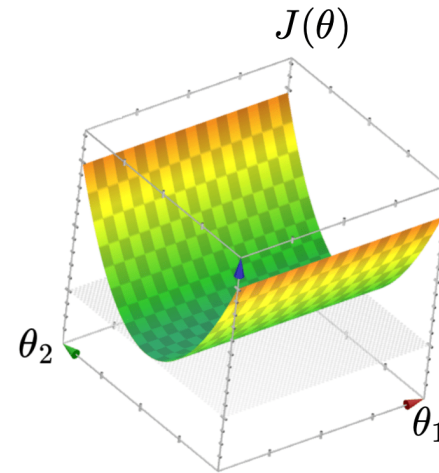
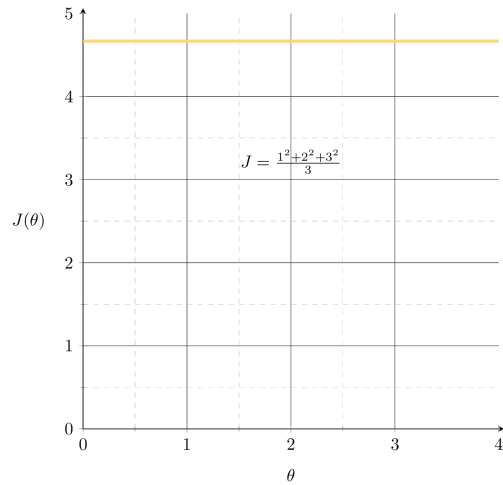
https://shenshen.mit.edu/demos/ridge/colinear_MSE.html

infinitely many optimal θ

data



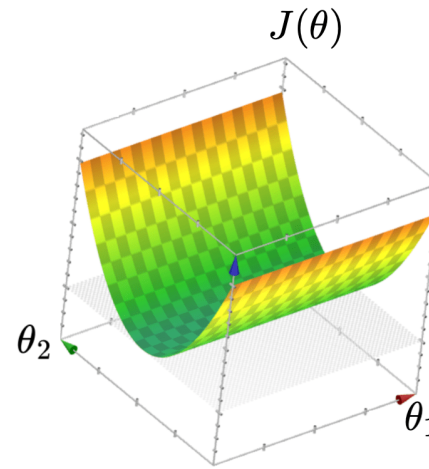
MSE



$$\theta^* = (X^\top X)^{-1} X^\top Y \text{ is not well-defined}$$

infinitely many optimal θ^*

Quick Summary:

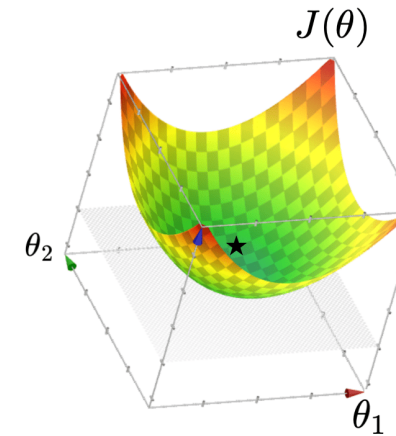


When X is not full column rank

- $J(\theta)$ has a "flat" bottom, like a half pipe
- 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
- $\theta^* = (X^\top X)^{-1} X^\top Y$
- θ^* gives the unique optimal hyperplane



$X^\top X$ becoming more invertible

when $X^\top X$ is *almost* singular, technically

$\theta^* = (X^\top X)^{-1} X^\top Y$ does exist

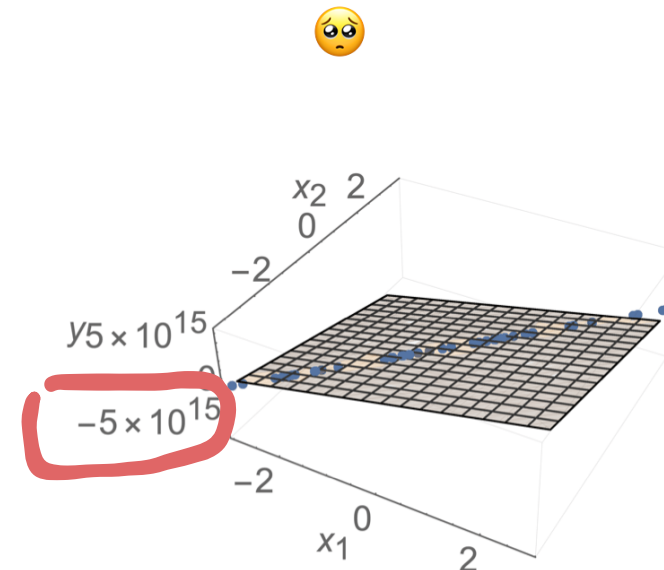
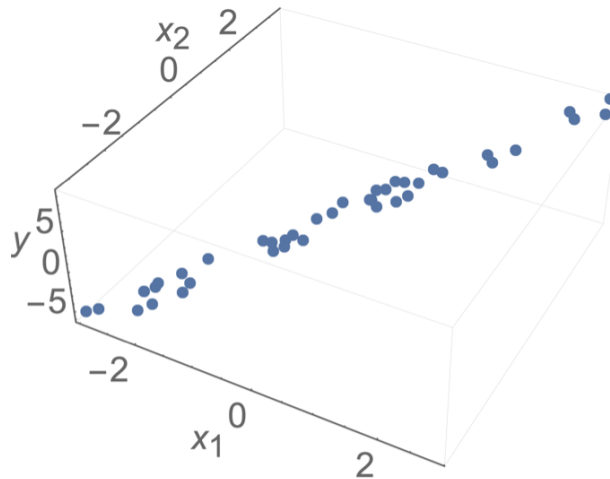
θ^* does give the unique optimal hyperplane

but

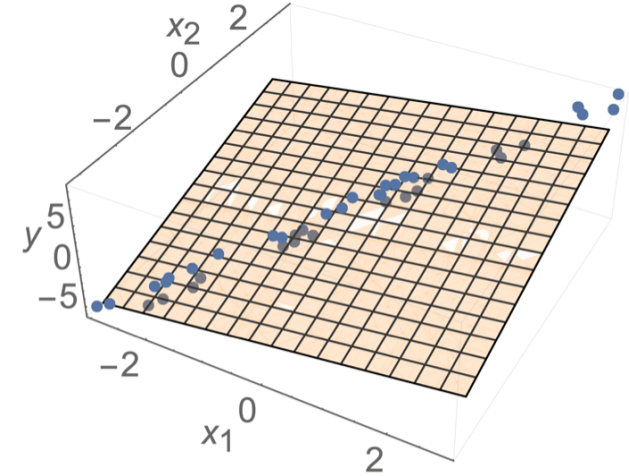
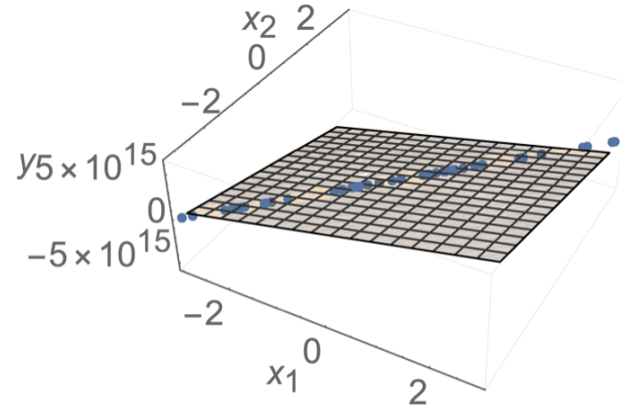
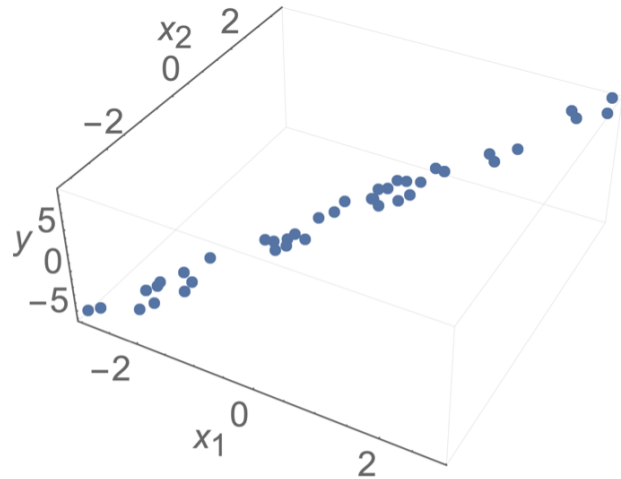
θ^* tends to be very sensitive to the small changes in the data

θ^* tends to have huge magnitude

θ^* tends to overfit



when $X^\top X$ is *almost* singular



lots of hypotheses (lots of θ s) fit the training data reasonably well

prefer θ with small magnitude (less sensitive prediction when x changes slightly)

Outline

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

Regularization

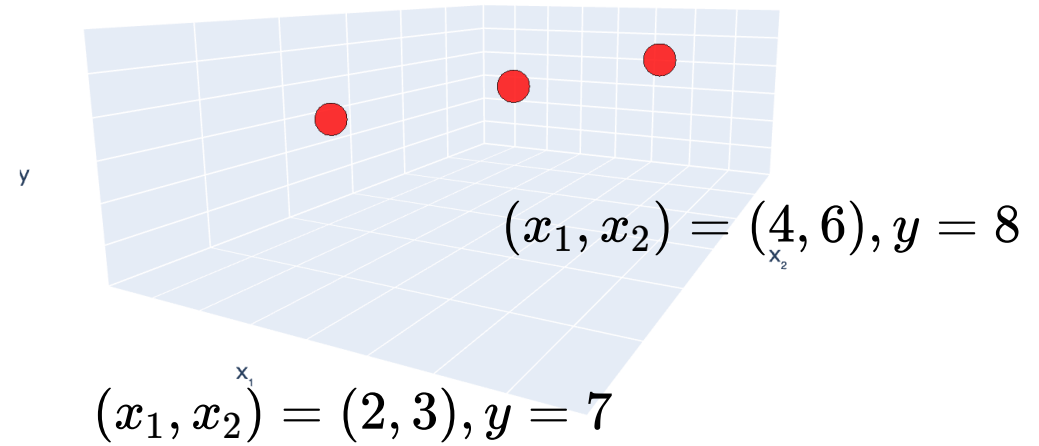
- technique to combat overfitting
- at a high-level, it's to sacrifice some training performance, in the hope that testing behaves better
- many ways to regularize (e.g. implicit regularization, drop-out)
- we will look at a particularly simple regularization today, the so-called ridge or l_2 -regularization

Ridge Regression

- Add a square penalty on the magnitude of the parameters
- $J_{\text{ridge}}(\theta) = \frac{1}{n}(X\theta - Y)^\top(X\theta - Y) + \lambda\|\theta\|^2$ $(\lambda > 0)$
- λ is a so-called "hyperparameter" (we've already seen a hyperparameter in lab 1)
- Setting $\nabla_{\theta} J_{\text{ridge}}(\theta) = 0$ we get $\theta_{\text{ridge}}^* = (X^\top X + n\lambda I)^{-1} X^\top Y$
- θ_{ridge}^* always exists, and is always the unique optimal parameters.
- (see ex/lab/hw for discussion about the offset.)

case (c) training data set again

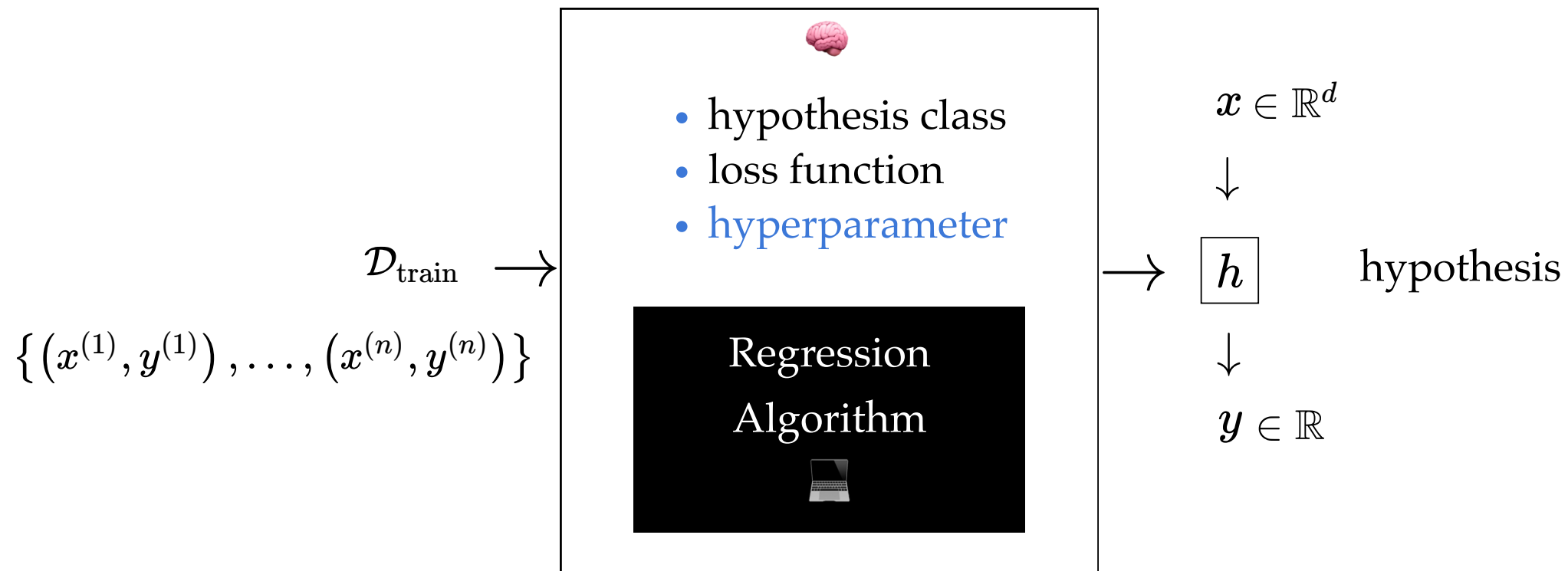
$$(x_1, x_2) = (6, 9), y = 9$$



<https://shenshen.mit.edu>

Comments on λ

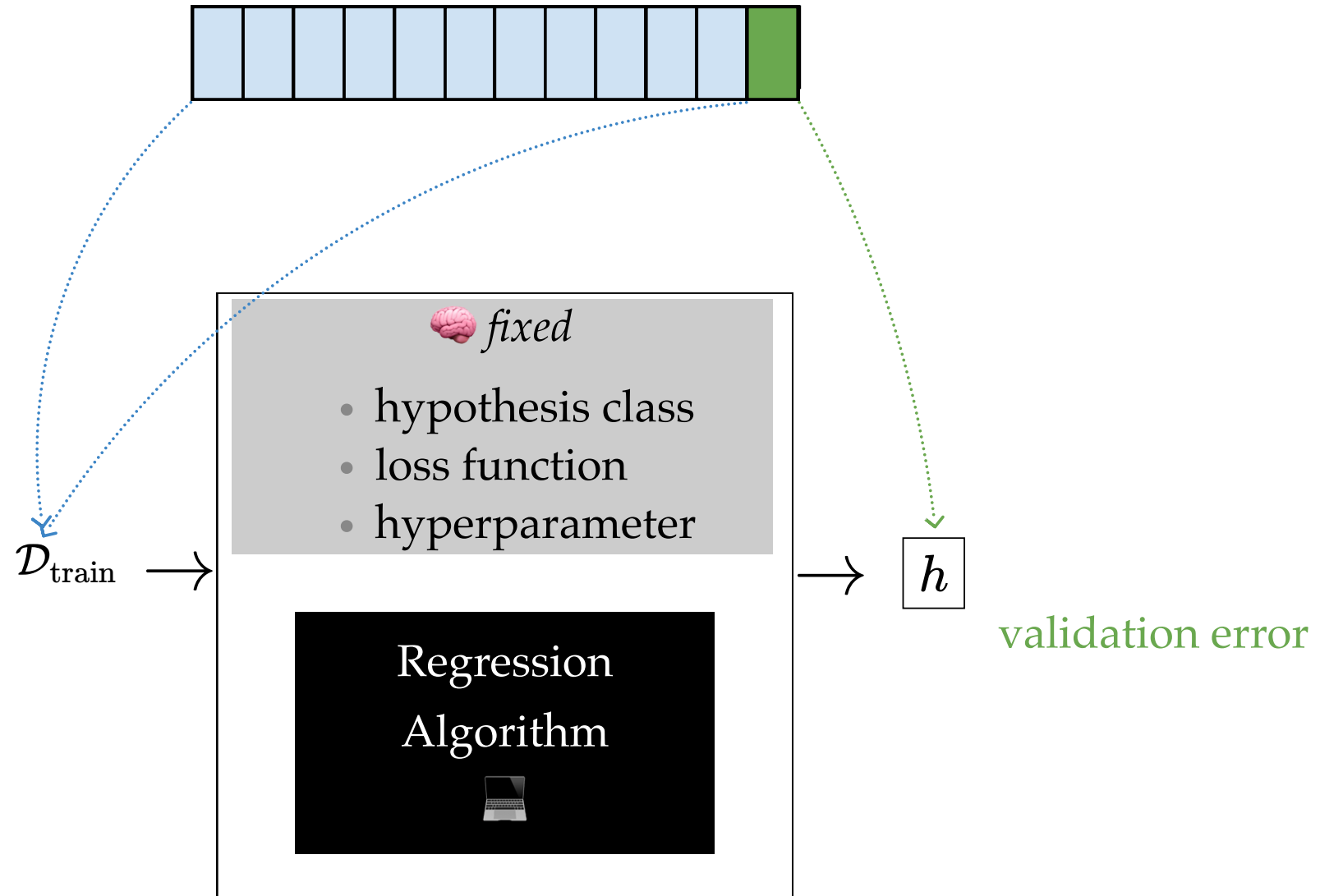
- one that's chosen by users, before we even see the data
- controls the tradeoff between MSE and theta magnitude
- implicitly controls the "richness" of the hypothesis class



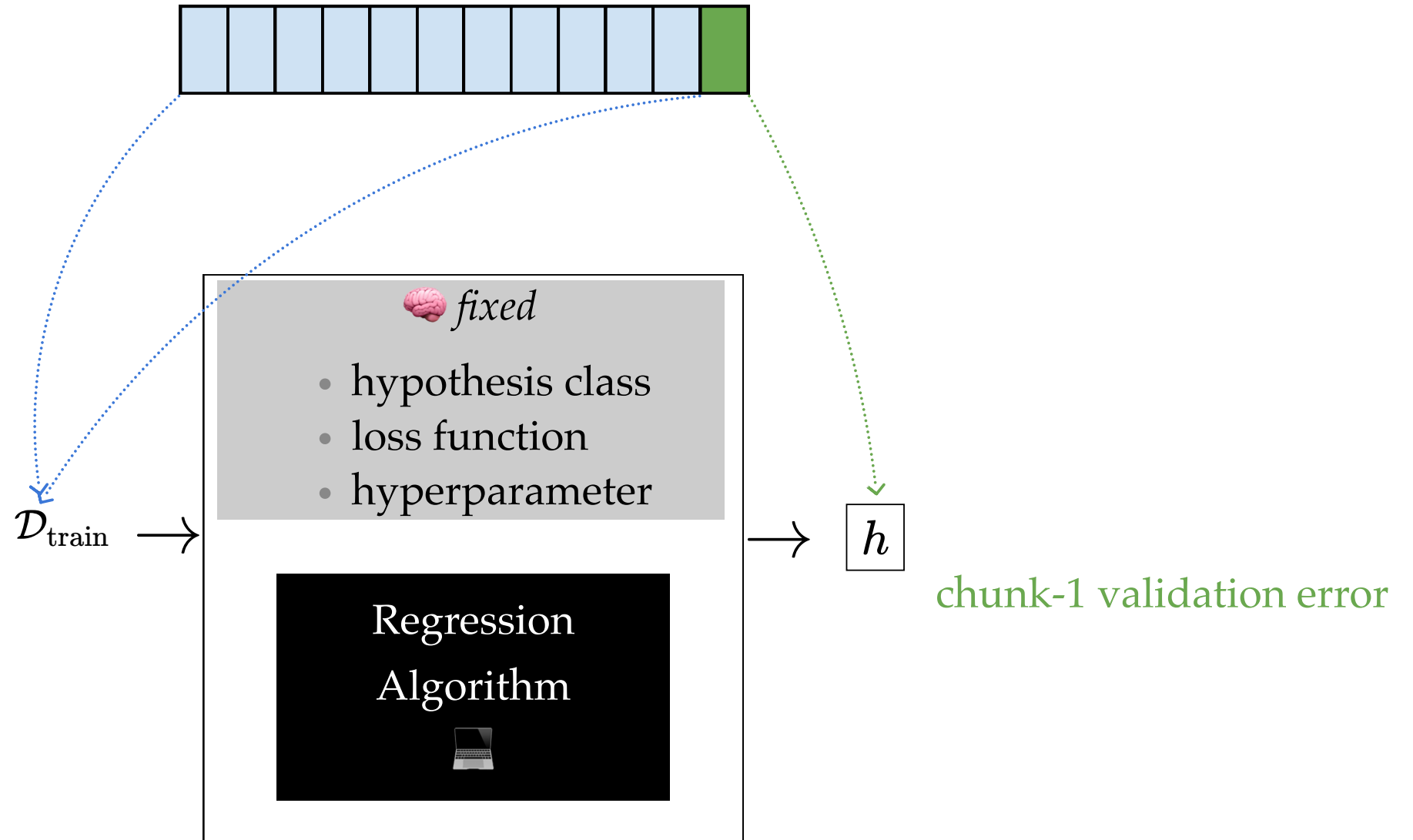
Outline

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

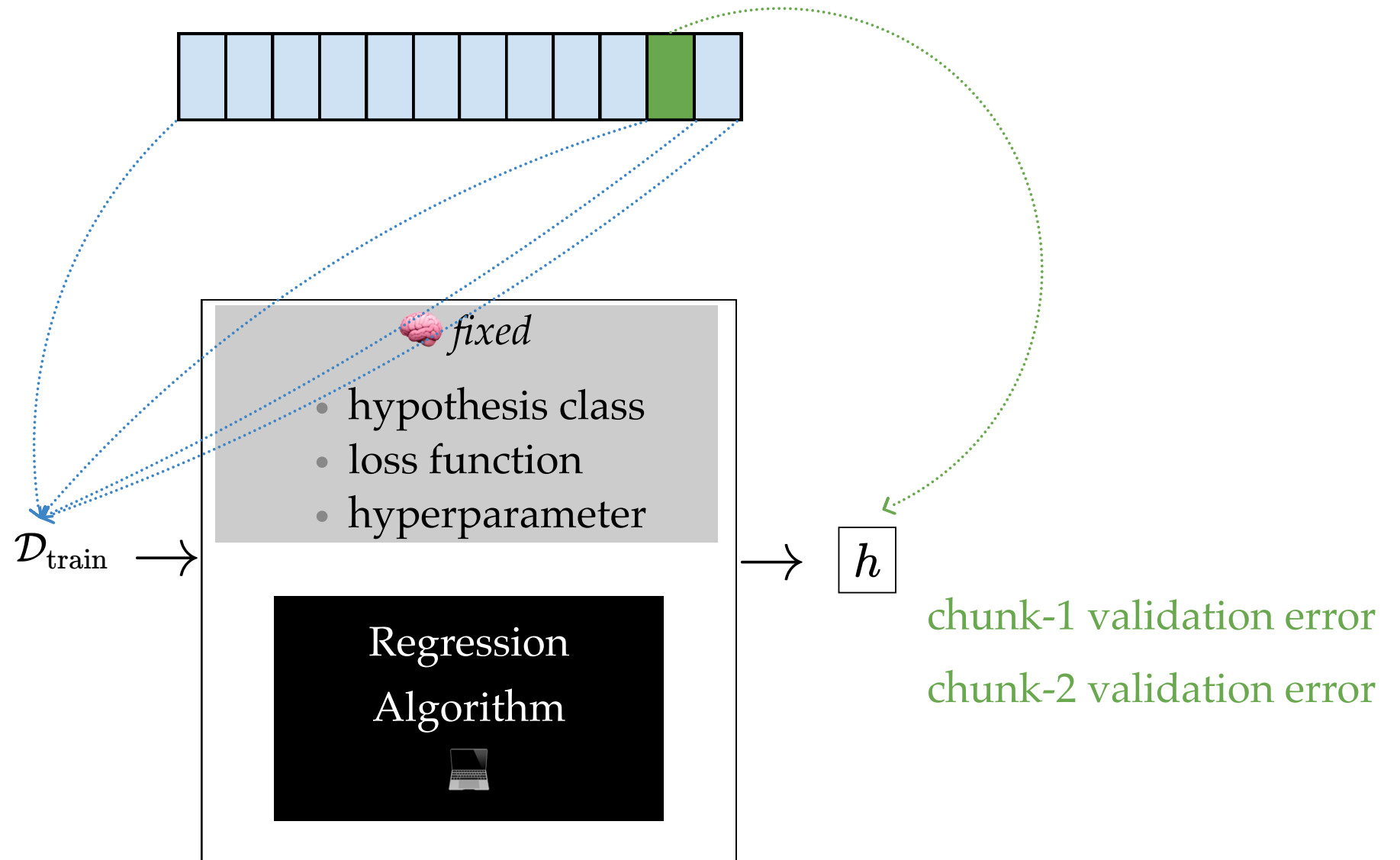
Validation



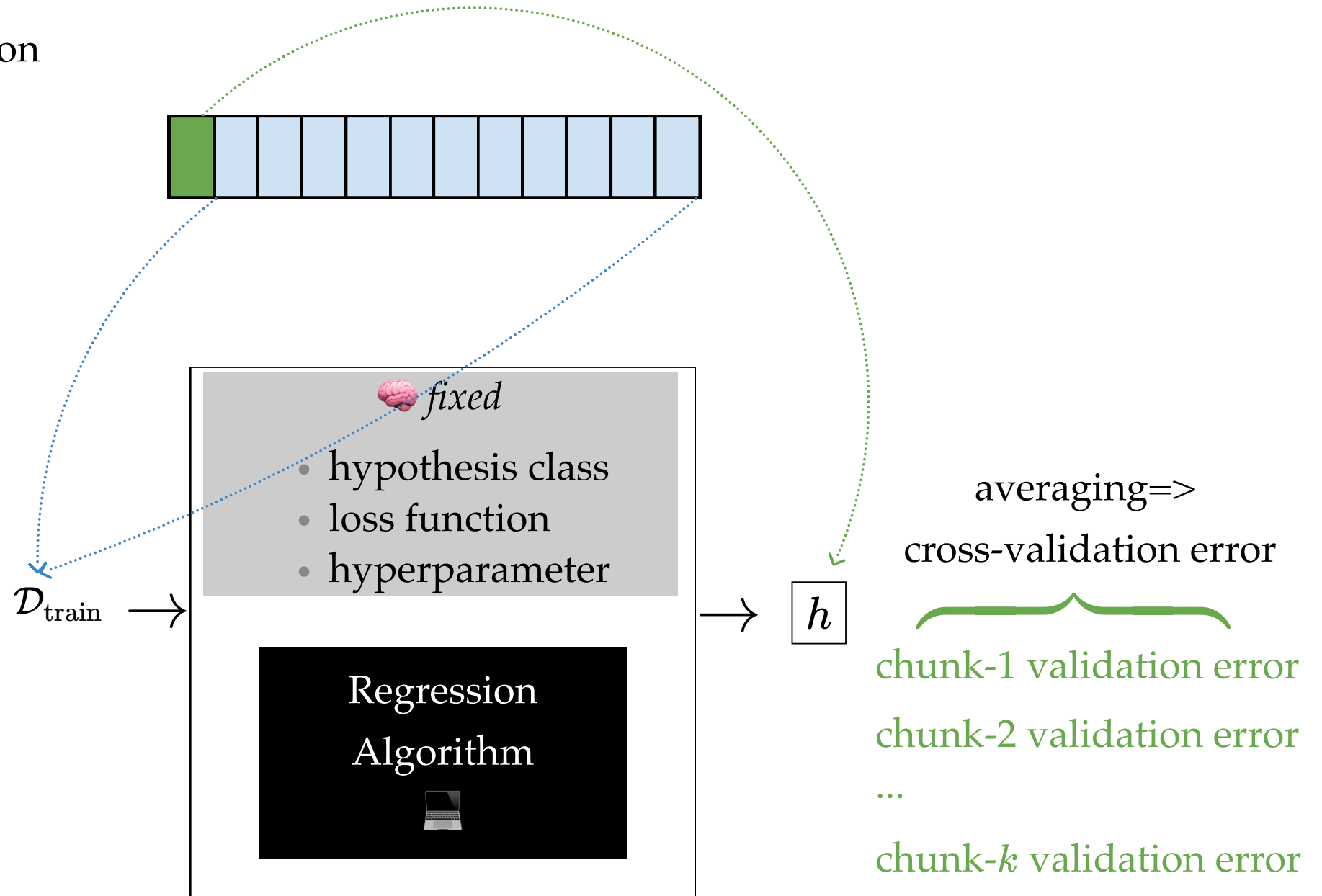
Cross-validation



Cross-validation



Cross-validation



Comments on cross-validation

- good idea to shuffle data first
- a way to "reuse" data
- cross-validation is more "reliable" than validation (less sensitive to chance)
- it's not to evaluate a hypothesis (testing error is)
- rather, it's to *evaluate* learning algorithm (e.g. hypothesis class choice, hyperparameter choice)
- Can have an outer loop for *picking* good hyperparameter or hypothesis class

Summary

- Closed-form formula for OLS is not well-defined when $X^T X$ is singular, and we have infinitely many optimal θ^* .
- Even in scenarios where $X^T X$ is just ill-conditioned, we get sensitivity issues, many almost-as-good solutions, while the absolutely best θ^* is overfitting to the data.
- We need to indicate our preference somehow, and also fight overfitting.
- Regularization helps battle overfitting -- by constructing a new optimization problem that implicitly prefers small-magnitude θ .
- Least-squares regularization leads to the ridge-regression formulation. (Good news: we can still solve it analytically!)
- λ trades off training MSE and regularization strength, it's a hyperparameter.
- Validation / cross-validation are a way to choose (regularization) hyperparameters.

https://docs.google.com/forms/d/e/1FAIpQLSftMB5hSccgAbIAFmP_LuZt95w6KFx0x_R3uuzBP8WwjSzZeQ/viewform?embedded=true

We'd love to hear
your **thoughts**.

Thanks!