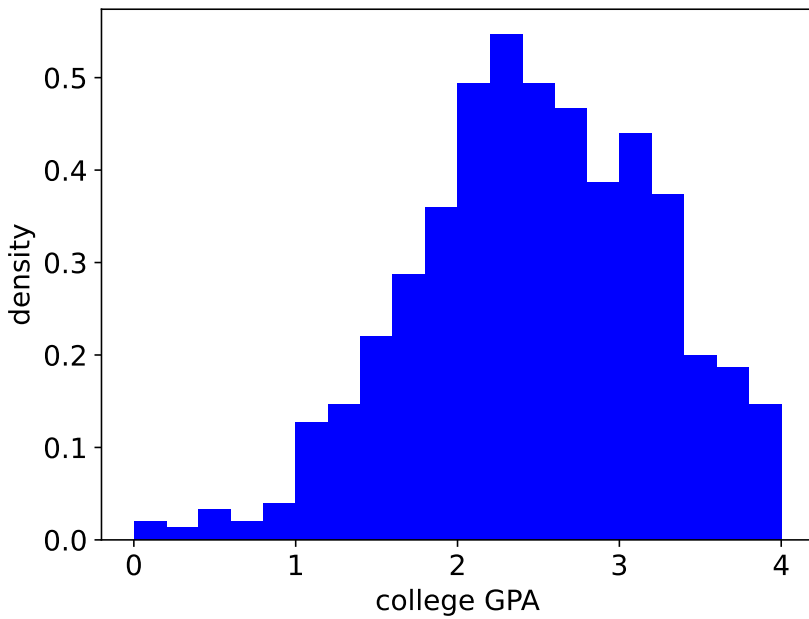


# Linear regression

COMS 4771 Fall 2023

**Dartmouth student dataset**

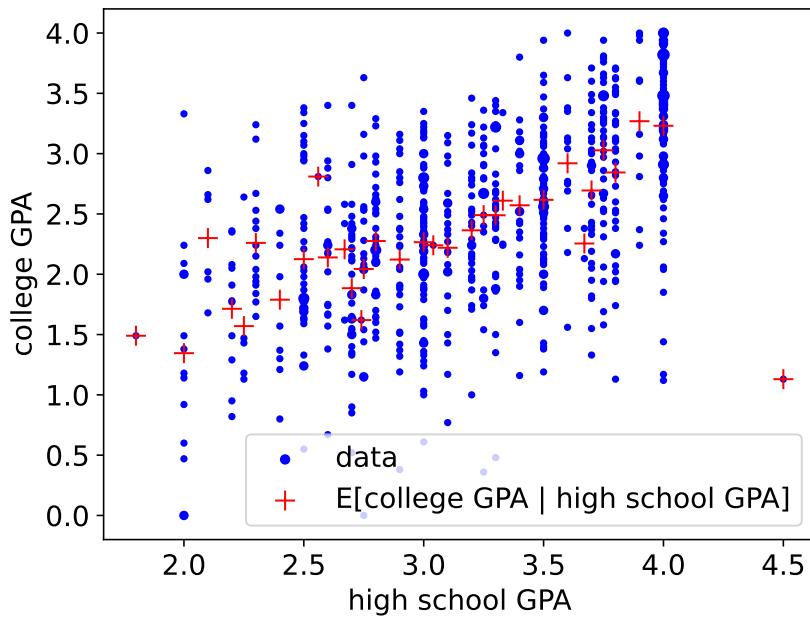
# Dataset of 750 Dartmouth students' (first-year) college GPA<sup>1</sup>



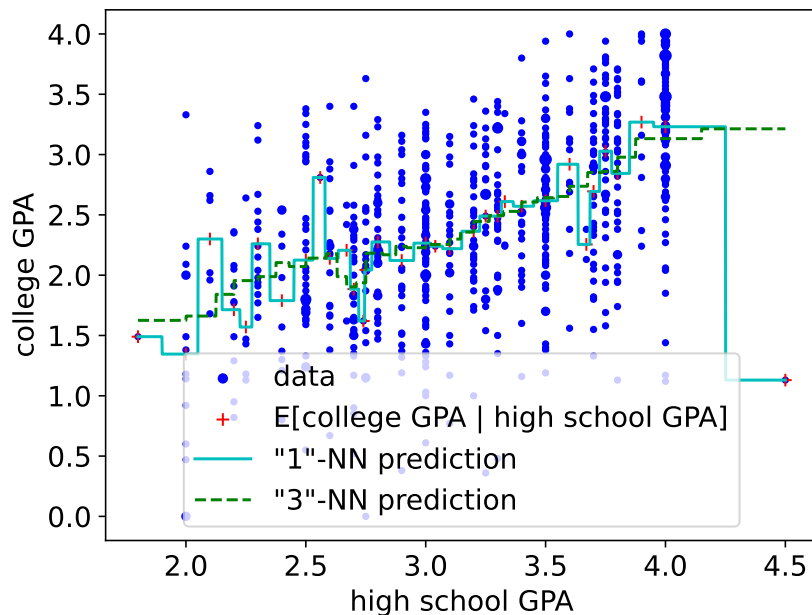
Mean 2.47  
Standard deviation 0.75

<sup>1</sup><https://chance.dartmouth.edu/course/Syllabi/Princeton96/ETSValidation.html>

Dartmouth dataset also has high school GPA of each student  
Question: Is high school GPA predictive of college GPA?



Attempting to exploit “local regularity” using NN



3 / 30

Possible “global” modeling assumption:

- ▶ Increase in high school GPA by  $\Delta$  should give an increase in (expected) college GPA by  $\propto \Delta$
- ▶ In other words,

$$\mathbb{E}[\text{college GPA} \mid \text{high school GPA}]$$

is \_\_\_\_\_ function of high school GPA

4 / 30

## Least squares linear regression

$f: \mathbb{R} \rightarrow \mathbb{R}$  is linear if it is of the form

$$f(x) = mx + b$$

for some parameters  $m, b \in \mathbb{R}$

Problem: given a dataset  $\mathcal{S}$  from  $\mathbb{R} \times \mathbb{R}$ , find (parameters of) a linear function  $f(x) = mx + b$  of minimal [sum of squared errors \(SSE\)](#)

$$\text{sse}[m, b] = \sum_{(x,y) \in \mathcal{S}} (mx + b - y)^2$$

Method of solution is called [ordinary least squares \(OLS\)](#)

6 / 30

Minimizers of SSE must be zeros of the two partial derivative functions:

$$\frac{\partial \text{sse}}{\partial m}[m, b] = 2 \sum_{(x,y) \in \mathcal{S}} (mx + b - y)x = 0$$

$$\frac{\partial \text{sse}}{\partial b}[m, b] = 2 \sum_{(x,y) \in \mathcal{S}} (mx + b - y) = 0$$

Two linear equations in two unknowns

Together, the equations are called the [normal equations](#)

7 / 30

Equivalent form:

$$\begin{aligned} \text{avg}(x^2) m + \text{avg}(x) b &= \text{avg}(xy) \\ \text{avg}(x) m + b &= \text{avg}(y) \end{aligned}$$

where

$$\begin{aligned} \text{avg}(x) &= \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} x, & \text{avg}(x^2) &= \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} x^2, \\ \text{avg}(xy) &= \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} xy, & \text{avg}(y) &= \frac{1}{|\mathcal{S}|} \sum_{(x,y) \in \mathcal{S}} y \end{aligned}$$

Solution to normal equations:

$$\begin{aligned} m &= \frac{\text{avg}(xy) - \text{avg}(x) \cdot \text{avg}(y)}{\text{avg}(x^2) - \text{avg}(x)^2}, \\ b &= \text{avg}(y) - m \cdot \text{avg}(x) \end{aligned}$$

What if  $\text{avg}(x^2) = \text{avg}(x)^2$ ?

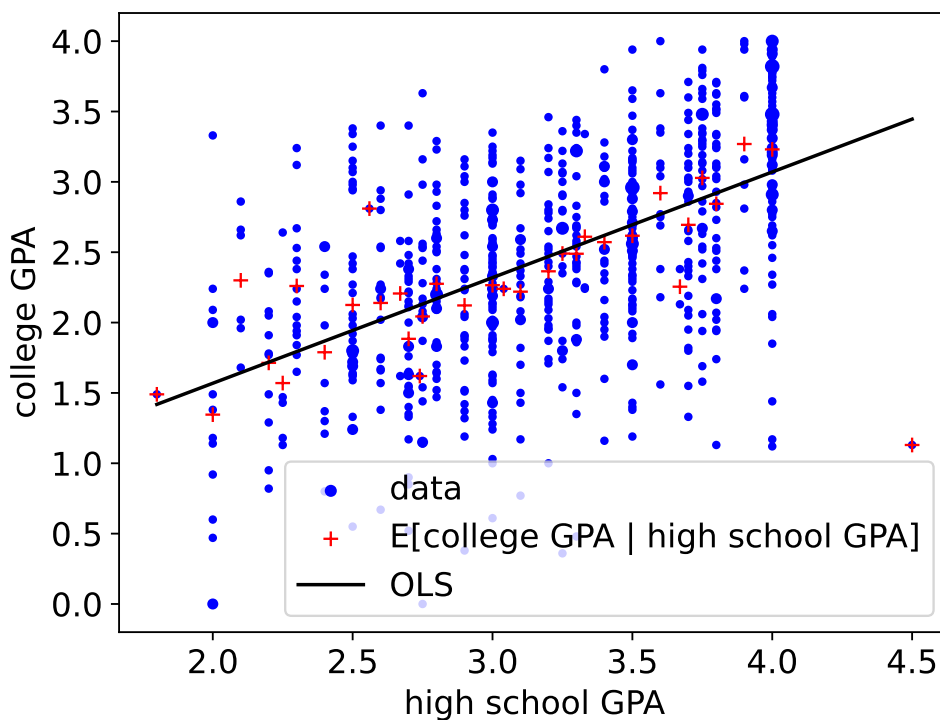
For Dartmouth dataset:

$$m = 0.751, \quad b = 0.067$$

RMSE:

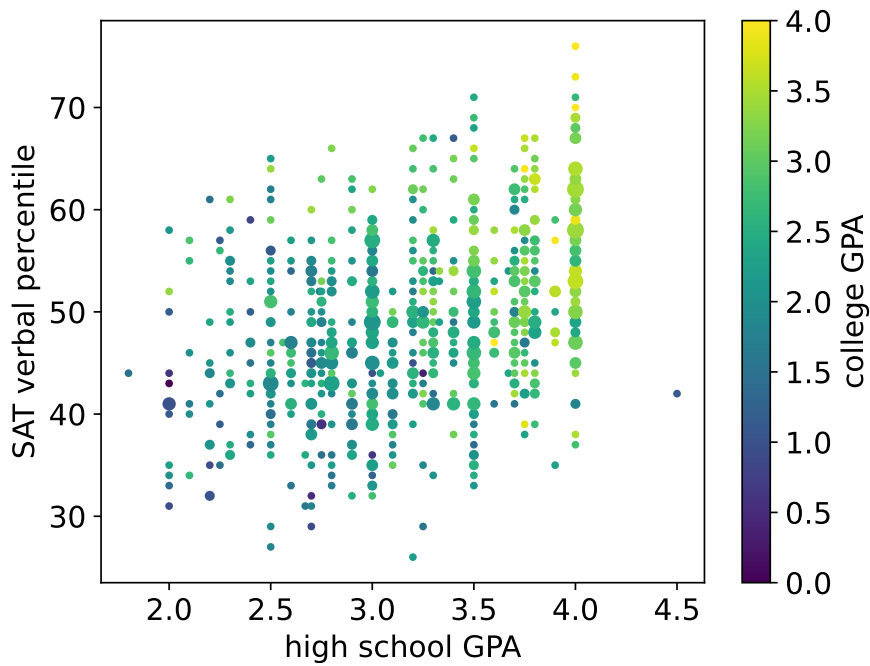
$$\sqrt{\frac{1}{|\mathcal{S}|} \text{sse}[m, b; \mathcal{S}]} = 0.629$$

(Recall standard deviation of college GPA is 0.75)



## Bivariate linear regression

Dartmouth dataset also includes SAT verbal percentiles





Linear function of two variables  $x_1$  and  $x_2$ :

$$f(x_1, x_2) = m_1x_1 + m_2x_2 + b$$

Problem: given a dataset  $\mathcal{S}$  from  $\mathbb{R}^2 \times \mathbb{R}$ , find (parameters of) a linear function  $f(x_1, x_2) = m_1x_1 + m_2x_2 + b$  of minimal sum of squared errors

$$\text{sse}[m, b; \mathcal{S}] = \sum_{(x_1, x_2, y) \in \mathcal{S}} (m_1x_1 + m_2x_2 + b - y)^2$$

13 / 30

Normal equations: three linear equations in three unknowns  $(m_1, m_2, b)$

$$\begin{bmatrix} \text{avg}(x_1^2) & \text{avg}(x_1x_2) & \text{avg}(x_1) \\ \text{avg}(x_2x_1) & \text{avg}(x_2^2) & \text{avg}(x_2) \\ \text{avg}(x_1) & \text{avg}(x_2) & 1 \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ b \end{bmatrix} = \begin{bmatrix} \text{avg}(x_1y) \\ \text{avg}(x_2y) \\ \text{avg}(y) \end{bmatrix}$$

Solve using elimination algorithm

14 / 30

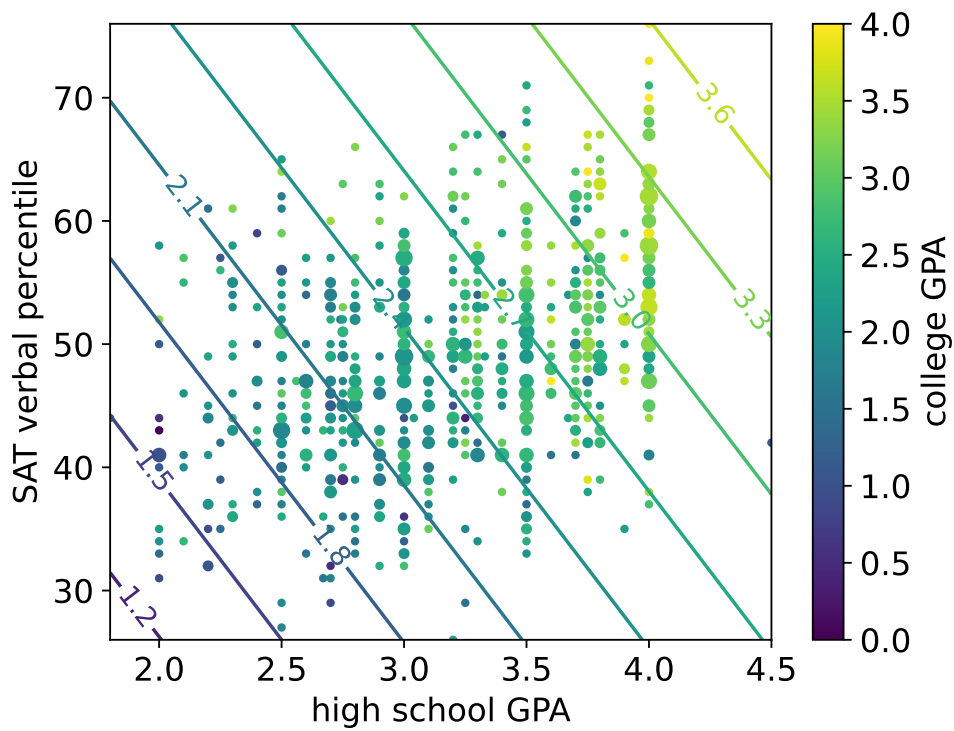
Dartmouth dataset:  $x_1 =$  high school GPA,  $x_2 =$  SAT verbal percentile

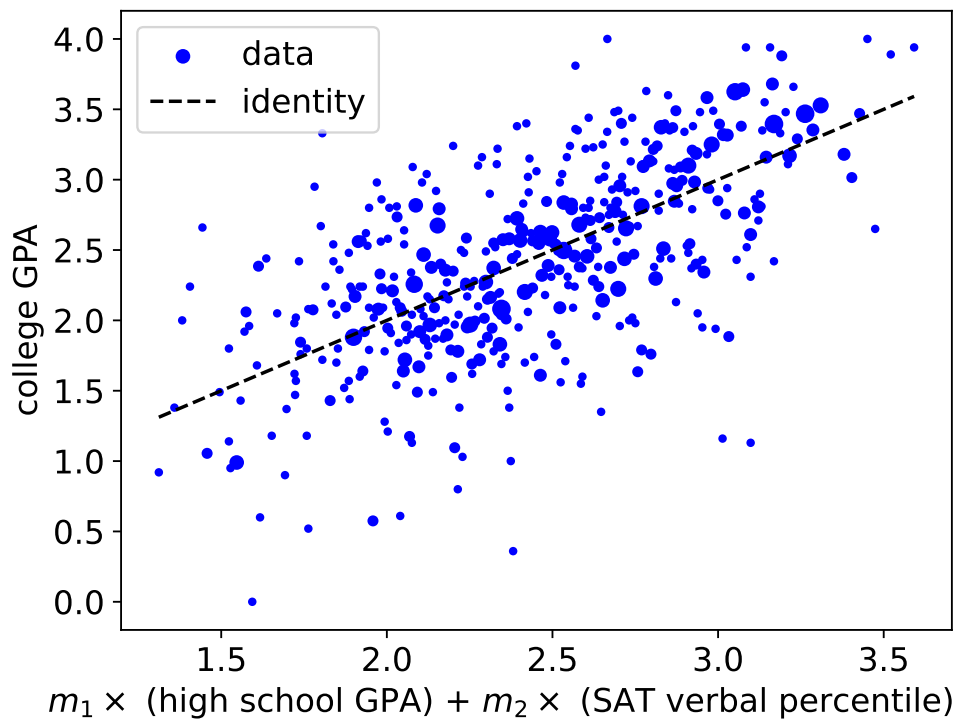
$$m_1 = 0.611, \quad m_2 = 0.024, \quad b = -0.639$$

RMSE:

$$\sqrt{\frac{1}{|\mathcal{S}|} \text{sse}[m_1, m_2, b; \mathcal{S}]} = 0.603$$

(Recall standard deviation of college GPA is 0.75)





## Linear algebra of ordinary least squares

(Homogeneous) linear function of  $d$  variables  $x = (x_1, \dots, x_d)$  is parameterized by  $d$ -dimensional weight vector  $w = (w_1, \dots, w_d)$ :

$$f_w(x) = w^\top x$$

To handle inhomogeneous linear functions (i.e., affine functions), include an extra always-1 feature:  $x_{d+1} = 1$

$$\begin{aligned} f_w(x) &= w^\top x \\ &= (w_1 x_1 + \dots + w_d x_d) + \underline{\hspace{2cm}} \end{aligned}$$

18 / 30

Problem: given a dataset  $\mathcal{S}$  from  $\mathbb{R}^d \times \mathbb{R}$ , find  $w \in \mathbb{R}^d$  of minimal sum of squared errors

$$\text{sse}[w; \mathcal{S}] = \sum_{(x,y) \in \mathcal{S}} (w^\top x - y)^2$$

Method of solution: OLS

19 / 30

**Matrix notation:** let  $\mathcal{S} = ((x^{(i)}, y^{(i)}))_{i=1}^n$ , and put

$$A = \begin{bmatrix} \leftarrow & (x^{(1)})^\top & \rightarrow \\ & \vdots & \\ \leftarrow & (x^{(n)})^\top & \rightarrow \end{bmatrix}, \quad b = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

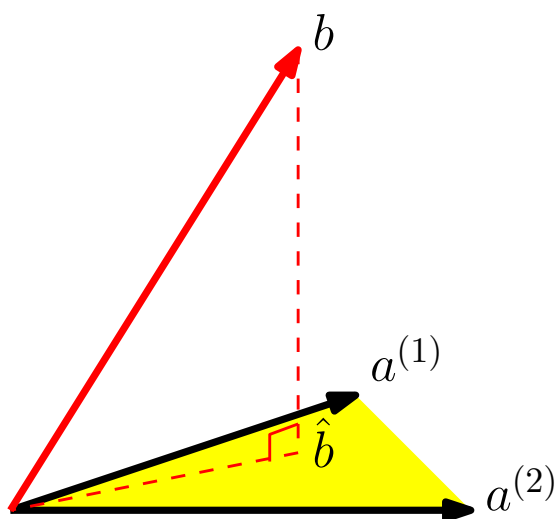
so

$$Aw = \begin{bmatrix} w^\top x^{(1)} \\ \vdots \\ w^\top x^{(n)} \end{bmatrix}, \quad Aw - b = \begin{bmatrix} w^\top x^{(1)} - y^{(1)} \\ \vdots \\ w^\top x^{(n)} - y^{(n)} \end{bmatrix}$$

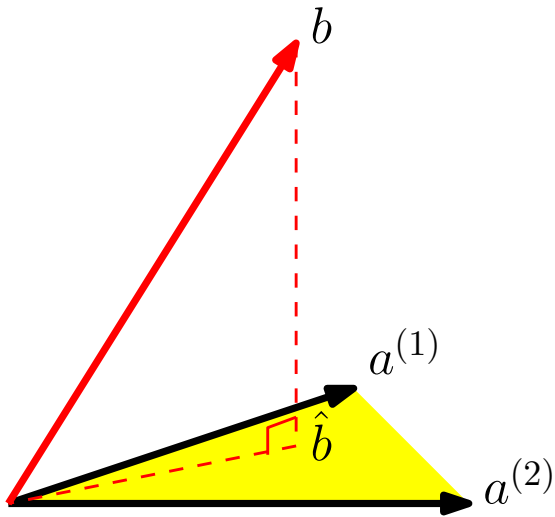
Therefore

$$\|Aw - b\|^2 = \sum_{i=1}^n \underline{\hspace{2cm}}$$

$Aw \in \text{CS}(A)$  for every  $w \in \mathbb{R}^d$



How many ways to write  $\hat{b}$  as a linear combination of the columns of  $A$ ?



22 / 30

### Normal equations in matrix notation

Key fact:  $\text{CS}(A)$  and  $\text{NS}(A^T)$  are orthogonal complements

23 / 30

Summary:

- ▶ Normal equations:  $(A^T A)w = A^T b$
- ▶ If  $\text{rank}(A) = d$ , then solution is unique
- ▶ Else, infinitely-many solutions
- ▶ Common choice for tie-breaking: minimum norm solution

$$\arg \min_{w \in \mathbb{R}^d} \|w\| \text{ s.t. } (A^T A)w = A^T b$$

24 / 30

```
def learn(train_x, train_y):  
    return np.linalg.pinv(train_x).dot(train_y)
```

```
def predict(params, test_x):  
    return test_x.dot(params)
```

25 / 30

## Statistical view of ordinary least squares

Normal linear regression model: Conditional distribution of  $Y$  given  $X = x$  is

$$N(w^\top x, \sigma^2)$$

- ▶  $w$  and  $\sigma^2$  are parameters of the model
- ▶ In this model, best possible MSE is  $\sigma^2$



## MLE in normal linear regression model

- ▶ Likelihood of  $w$  and  $\sigma^2$ :

$$L(w, \sigma^2) = \prod_{(x,y) \in \mathcal{S}} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - w^\top x)^2}{2\sigma^2}\right)$$

- ▶ Log-likelihood:

$$\ln L(w, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{(x,y) \in \mathcal{S}} (y - w^\top x)^2 - \frac{|\mathcal{S}|}{2} \ln(2\pi\sigma^2)$$

- ▶ In terms of  $w$ , maximizing log-likelihood is same as minimizing SSE!

27 / 30

## Statistical inference (example)

- ▶ Suppose you fit linear regression model to data, and find that  $w \neq (0, \dots, 0)$

How confident are you in this finding?

28 / 30

## Generalization

- ▶ Suppose  $\mathcal{S} \stackrel{\text{i.i.d.}}{\sim} (X, Y)$
- ▶ OLS gives minimizer of empirical risk (for square loss, among linear functions)

$$\widehat{\text{Risk}}[w] = \frac{1}{n} \sum_{(x,y) \in \mathcal{S}} \text{loss}_{\text{sq}}(w^\top x, y)$$

But we actually care about the (true) risk

$$\text{Risk}[w] = \mathbb{E}[\text{loss}_{\text{sq}}(w^\top X, Y)]$$

- ▶ Is empirical risk a good estimate of (true) risk?
  - ▶ Usually only if  $|\mathcal{S}|$  is sufficiently large

**Extreme example:**  $d = 1$ ,  $|\mathcal{S}| = 2$ ,  $\widehat{\text{Risk}}[w] = 0$

