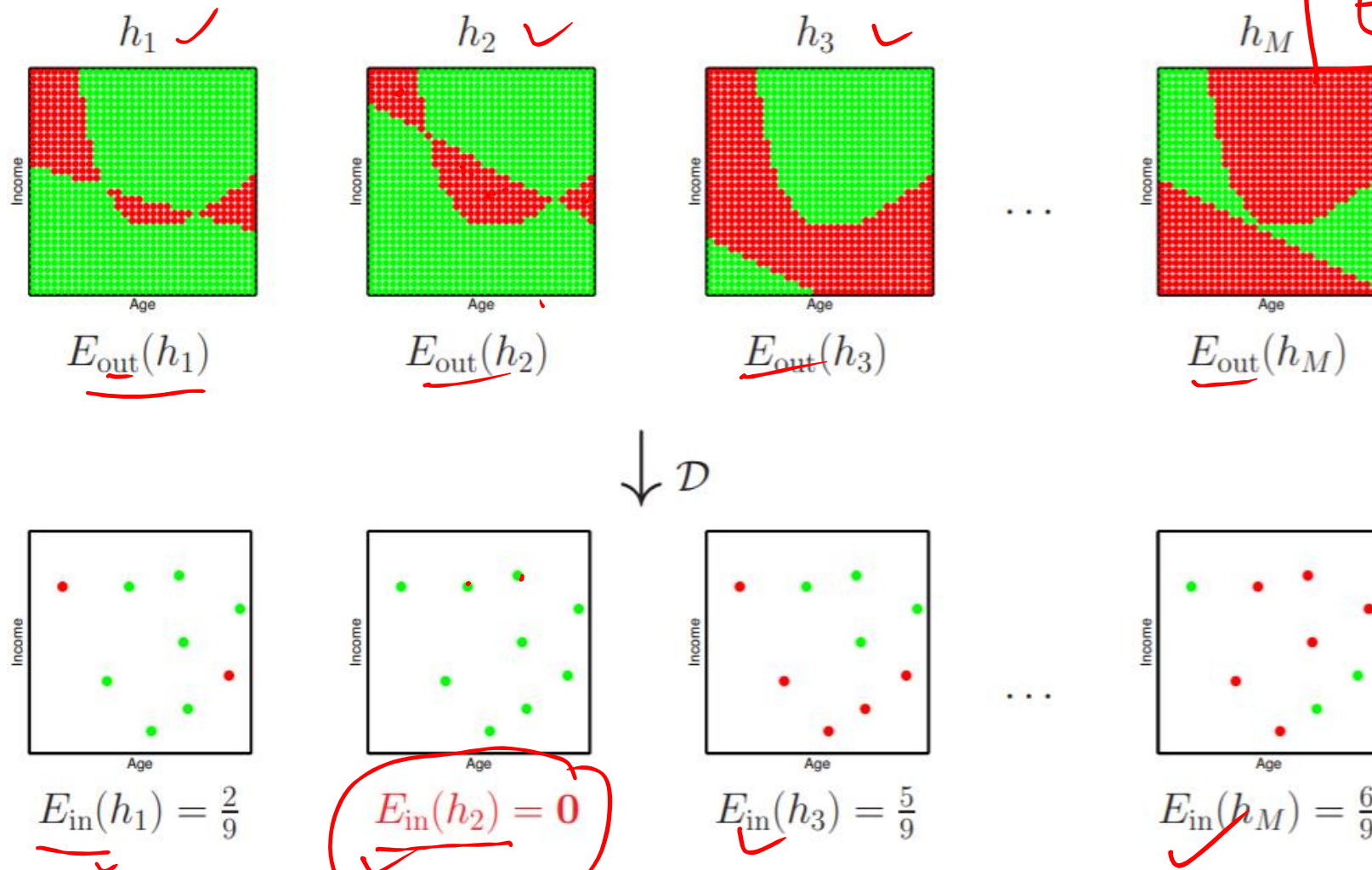


Machine Learning from Data

Lecture 4: Spring 2021

Today's Lecture

- Feasibility of Learning ✓
- Two Step Solution to Learning
- Error and Noise



Pick the hypothesis with minimum E_{in} ; will E_{out} be small?

Verification and Selection Bias

- If we pick the hypothesis with minimum in-sample error, it does not approximate out-of-sample error.
- Search Causes Selection Bias
- In Real Learning in-sample error cannot reach out to out-of-sample error.

$$E_{in}(h_2) = 0$$

$$P(\text{Bad}) \leq \text{something small}$$

$$\text{Prob} [|E_{out}(g) - E_{in}(g)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

Not Valid

Using Hoeffding's Inequality in Learning

- Definition - "Hoeffding's inequality provides an upper bound on the probability that the sum of bounded independent random variables deviates from its expected value by more than a certain amount."

{ 1) For the validity of the bound we must
fix our hypothesis before we see the data
2) To change h after looking at data
Violating Hoeffding's inequality.

Updating Hoeffding's Bound

• Set $\mathcal{H} \xrightarrow{\text{LA}}$ picks g based on D after generating data.
 $\times P[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon]$ is small for any fixed h_m

$$P[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon]$$

Bound is updated

is small for the final hypothesis g .

$$\mathcal{H}, h_m \in \mathcal{H}$$

$$\underbrace{|E_{in}(g) - E_{out}(g)| > \epsilon}_{\text{subset relation}} \Rightarrow \begin{aligned} &|E_{in}(h_1) - E_{out}(h_1)| > \epsilon \\ \text{or} &|E_{in}(h_2) - E_{out}(h_2)| > \epsilon \\ &\vdots \\ \text{or} &|E_{in}(h_m) - E_{out}(h_m)| > \epsilon \end{aligned} \quad \text{--- ①}$$

$A, B \rightarrow \text{events}$

* If $A \Rightarrow B$ or $A \subseteq B$, $P[A] \leq P[B]$

* A, B, C, \dots, Z : $P(A \text{ or } B \text{ or } C \dots \text{ or } Z) \leq P[A] + P[B] + \dots + P[Z]$

= Union bound.

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \stackrel{(1)}{\leq} P[|E_{in}(h_1) - E_{out}(h_1)| > \epsilon] + P[|E_{in}(h_2) - E_{out}(h_2)| > \epsilon] + \dots + P[|E_{in}(h_M) - E_{out}(h_M)| > \epsilon]$$

$$\leq \sum_{m=1}^M P[|E_{in}(h_m) - E_{out}(h_m)| > \epsilon] \leq 2M e^{-2\epsilon^2 N} \quad (2)$$

Feasibility of Learning

- Two Questions to answer:

- • Can we make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$ $E_{out}(g) \approx E_{in}(g)$
- • Can we make $E_{in}(g)$ small enough ✓

Complexity of H and f

$|H|$

complex phenomena

Bigger choice set

Price

Complexity of K

$|H| \uparrow$

→ Yes → Q_1 , keep M in check.

→ Yes → Q_2 , $M \uparrow$ due to the result of choice

Complexity of f

How hard is f ?

→ Yes → Q_1 → No direct impact because not there in the off-diagonal inequality.

→ D is a result of f
 f is complex

$E_{in} \uparrow$

$|H| \uparrow \rightarrow$ bad for Q_1

Interpreting the Hoeffding's Bound

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2|\mathcal{H}|e^{-2\epsilon^2 N}, \quad \text{for any } \epsilon > 0.$$

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| \leq \epsilon] \geq 1 - 2|\mathcal{H}|e^{-2\epsilon^2 N}, \quad \text{for any } \epsilon > 0.$$

$$\delta = 2|\mathcal{H}|e^{-2\epsilon^2 N}$$

Theorem: With Probability at least $1 - \delta$

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{1}{2N} \log \frac{2|\mathcal{H}|}{\delta}}$$

Proof: $P[|E_{in}(g) - E_{out}(g)| \leq \epsilon] \geq 1 - \delta$
In other words: With probability at least $1 - \delta$

$$|E_{in}(g) - E_{out}(g)| \leq \epsilon$$

$$\Rightarrow E_{out}(g) \leq E_{in}(g) + \epsilon$$

$$\begin{aligned} \delta &= 2^{-N} \\ \epsilon &= \sqrt{\frac{1}{2N} \log \frac{2}{\delta}} \end{aligned}$$

$$E_{out}(g) \leq E_{in}(g) + \underbrace{\sqrt{\frac{1}{2N} \log \frac{2}{\delta}}}_{\text{error bar}}$$

E_{in} reaches out to E_{out} when H is small

$$\underline{E_{out}(g)} \leq \underline{E_{in}(g)} + \underbrace{\sqrt{\frac{1}{2N} \left(\log \frac{2|\mathcal{H}|}{\delta} \right)}}_{\text{error bar} \uparrow}$$

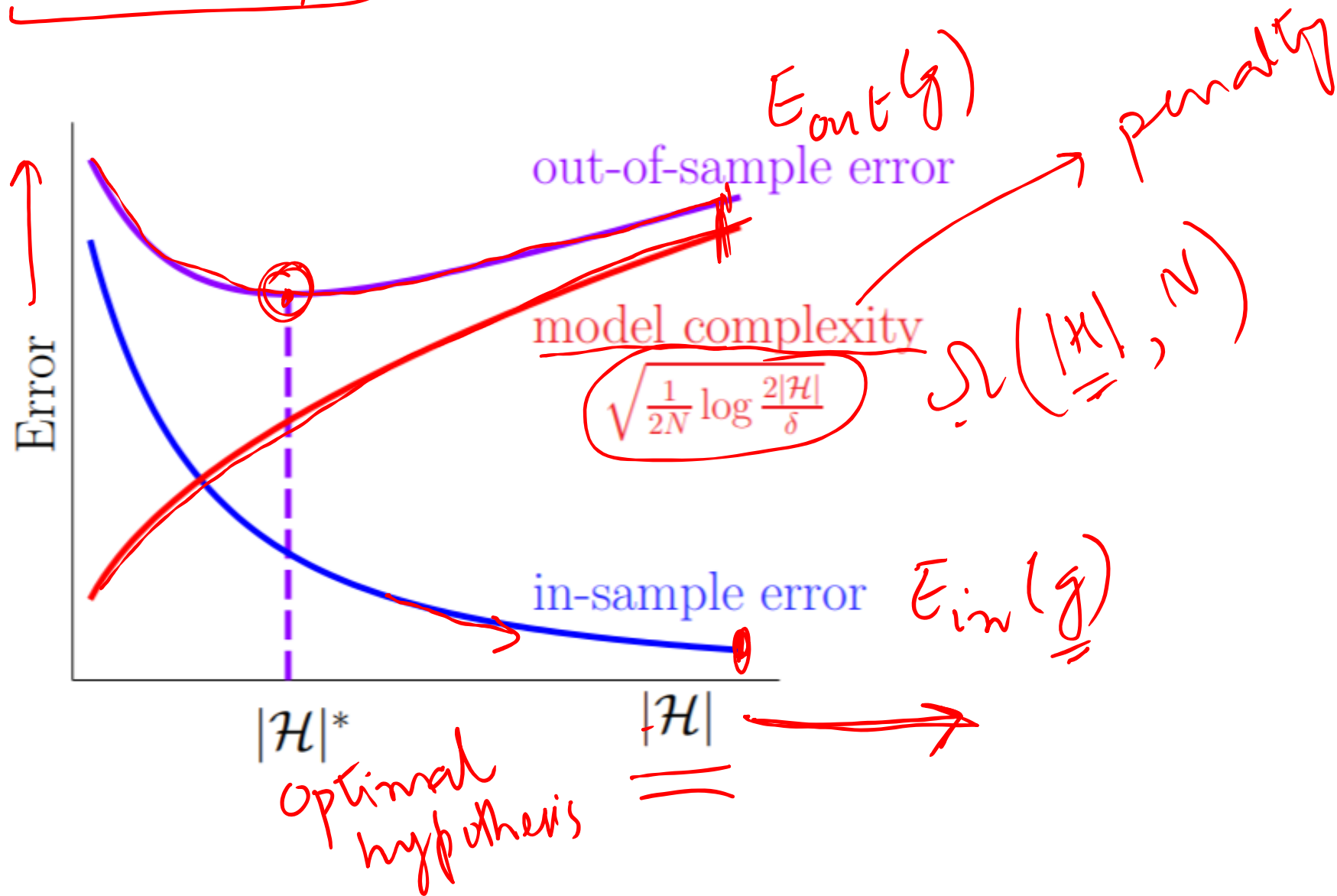
generalization bound.

If $N \gg \ln |\mathcal{H}|$

$$\underline{E_{out}(g)} \approx \underline{E_{in}(g)}$$

$E_{out}(g) ?$

Trade-off



2 Step Approach

- - 1) Ensure $E_{out}(g) \simeq E_{in}(g)$
 - 2) Get $E_{in}(g) \simeq 0$

Implies that

sample \rightarrow universe
 $E_{in}(g) \simeq E_{out}(g) \checkmark$

$g \simeq f$

$E_{out}(g) \simeq 0 ?$

$E_{out}(g) \simeq 0$

Summarize

- Is Learning Feasible? (Finite \mathcal{H})

(A) Yes, provide we accept a 2-step approach

$$\left\{ \begin{array}{l} \rightarrow E_{\text{out}}(g) \approx E_{\text{in}}(g) \\ \rightarrow E_{\text{in}}(g) \approx 0 \end{array} \right.$$

{ Theoretically & with high probability

$$\left\{ \begin{array}{l} \rightarrow E_{\text{in}}(g) \approx 0 \end{array} \right.$$

Conditions: 1) Fin \mathcal{H} ahead of time $\Rightarrow g \in \mathcal{H}$

2) Data must be IID ($P(x)$)

3) $E_{\text{out}}(g) \rightarrow x \rightarrow$ comes from $P(x)$

{ Sample Bin

B) \mathcal{M} has a possibility to fail

$A + B \longrightarrow$ Summarize the fact that
Learning is feasible.

Our Learning Approach is General

- 1) Applies to any target function f .
- 2) Applies to any $P(x)$
- 3) Applies to any hypothesis set \mathcal{H} (PLA, NN, SVM)
- 4) Applies to any LA.

All our conditions must hold.

Target Function

- 1) Complexity of target function \rightarrow Need more data
- 2) Noisy f. (stochastic) \rightarrow Practice
- 3) Error

Complex $f \rightarrow$ hard to learn, $E_{in}(g) \approx 0$

$|H| \uparrow \rightarrow \Omega(|H|, N) \uparrow$


Big N

$\sqrt{\frac{|H|}{N}} \rightarrow$ bigger

$E_{in} \approx E_{out}$

Noisy target

•

Same x 

$$f = P[y|x] \leftarrow$$

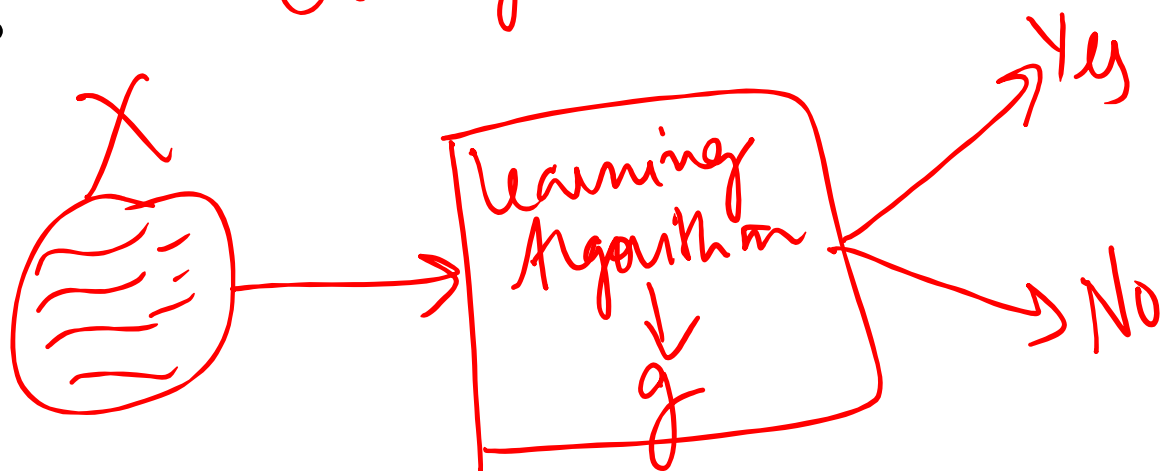
$$E_{in} \neq 0$$

$f \rightarrow$ noisy

Error

(User defined)

Classification problem



Supermarket (SM)

- 1) Say No to the right person
costly for the SM (False⁺_{ne})
- 2) Say Yes to the wrong person
not costly → okay (False_{true})

CIA access control

- 1) Say Yes to the wrong
person, extremely costly
for CIA (False_{true})⁺
- 2) Say No to the right
person, okay
(False_{ne})

Interpretation of Error (Not good)

Risk matrix (SM)

•

		g	
		Yes	No
Reality	Yes	0	10
	No	1	0

Risk Matrix (LIA)

		g	
		Yes	No
Reality	Yes	0	1
	No	1000	0

Pointwise Errors

$$e(h(x), y)$$

- 1) Binary errors (classification problems) $\left. \begin{array}{l} e(h(x), y) = [h(x) \neq y] \end{array} \right\} \text{indicator function}$

- 2) Squared errors (regression problems)
$$e(h(x), y) = (h(x) - y)^2$$

$$E_{\text{in}}(h) = \frac{1}{N} \sum_{i=1}^N e(h(x_i), y_i) \leftarrow$$

$$E_{\text{out}}(h) = E_x[e(h(x), f(x))]$$

Learning Set-up

2-Step Process

1) $E_{out}(g) \approx E_{in}(g) *$

2) $E_{in} \approx 0$

can ~~not~~ be a failure

