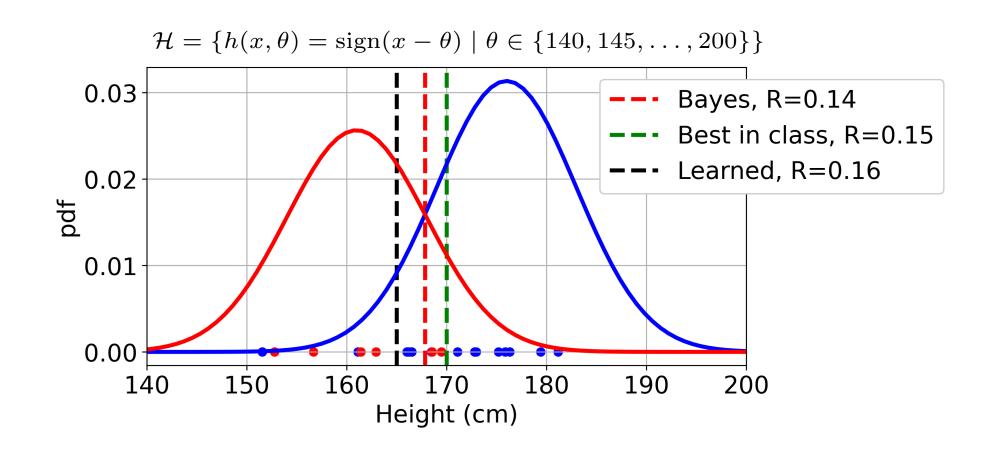
# Statistical Machine Learning (BE4M33SSU) Lecture 3: Probably Approximately Correct Learning

Czech Technical University in Prague V. Franc

## **Error decomposition**

#### **Errors:**

- 1. Best (Bayes) attainable risk  $R(p, h_*)$ , where  $h_*(x) = \underset{h \in \mathcal{Y}^{\mathcal{X}}}{\arg \min} R(p, h)$
- 2. Best risk in the class  $R(p, h_{\mathcal{H}})$ , where  $h_{\mathcal{H}} = \underset{h \in \mathcal{H}}{\arg \min} R(p, h)$
- 3. Risk of the learned predictor  $R(p,h_m)$ , where  $h_m=A(T_m)$



## **Error decomposition**

#### **Errors:**

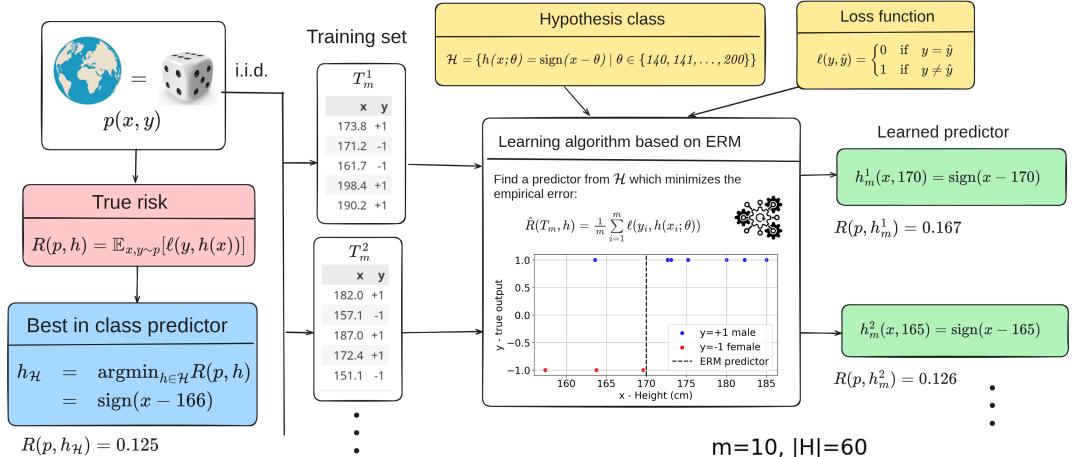
- 1. Best (Bayes) attainable risk  $R(p, h_*)$ , where  $h_*(x) = \underset{h \in \mathcal{Y}^{\mathcal{X}}}{\operatorname{arg \, min}} R(p, h)$
- 2. Best risk in the class  $R(p, h_{\mathcal{H}})$ , where  $h_{\mathcal{H}} = \underset{h \in \mathcal{H}}{\arg \min} R(p, h)$
- 3. Risk of the learned predictor  $R(p,h_m)$ , where  $h_m=A(T_m)$

#### **Error decomposition:**

$$\underbrace{R(p,h_m)}_{\text{learned predictor risk}} = \underbrace{\left(R(p,h_m) - R(p,h_{\mathcal{H}})\right)}_{\text{estimation error}} + \underbrace{\left(R(p,h_{\mathcal{H}}) - R(p,h_*)\right)}_{\text{approximation error}} + \underbrace{R(p,h_*)}_{\text{Bayes risk}}$$

- lacktriangle The approximation error: depends on  ${\cal H}$  chosen prior to learning.
- The estimation error: depends on  $\mathcal{H}$ , training data  $T_m$  and the algorithm A.
- Best (Bayes) attainable risk: irreducible error.

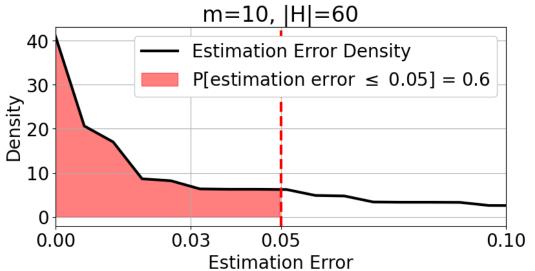
## **Probably Approximately Correct**



Given  $\varepsilon > 0$ , learned predictor  $h_m$  is approximately correct provided:

$$\underbrace{R(p, h_m) - R(p, h_{\mathcal{H}})}_{\text{estimation error}} \leq \varepsilon$$

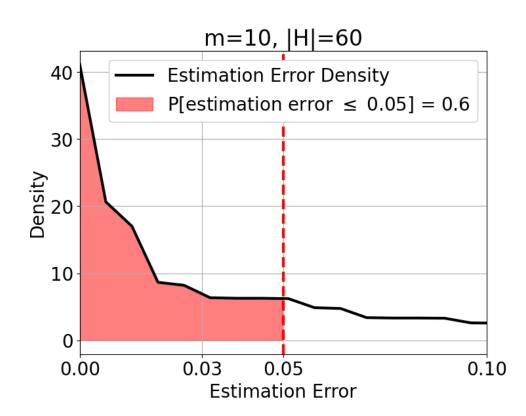
$$\underbrace{\text{estimation error}}_{\text{approximately correct}}$$

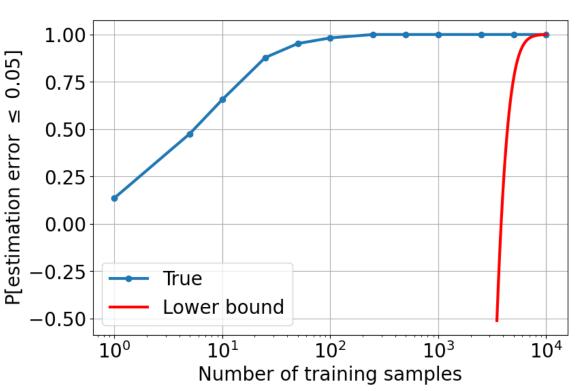


## **Probably Approximately Correct**



5/16





- ERM algorithm:  $h_m = A(T_m) = \operatorname*{arg\,min}_{h \in \mathcal{H}} \left[ \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq h(x_i) \rrbracket \right]$
- We derive a lower bound valid for any distribution p(x,y) and finite hypothesis class  $\mathcal{H}=\{h_1,\ldots,h_H\}$ :

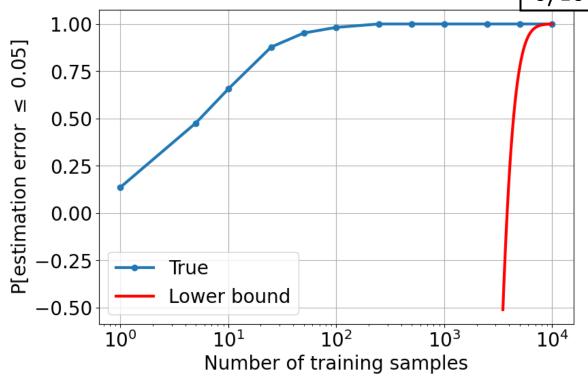
$$\mathbb{P}_{T_m \sim p^m} \left[ \underbrace{R(p, h_m) - R(p, h_{\mathcal{H}}) \leq \varepsilon}_{\text{approximately correct}} \right] \geq \underbrace{1 - 2 |\mathcal{H}| \, e^{-\frac{1}{2} m \, \varepsilon^2}}_{\text{lower bound increasing with } m}$$

#### Setup:

• 
$$h_m = \underset{h \in \mathcal{H}}{\operatorname{arg min}} \left[ \frac{1}{m} \sum_{i=1}^{m} [y_i \neq h(x_i)] \right]$$

Finite hypothesis class:

$$\mathcal{H} = \left\{ h^i \colon \mathcal{X} \to \mathcal{Y} \mid i \in \{1, \dots, H\} \right\}$$



Distribution-free lower bound:

$$\mathbb{P}_{T_m \sim p^m} \left[ \underbrace{R(p, h_m) - R(p, h_{\mathcal{H}}) \leq \varepsilon}_{\text{approximately correct}} \right] \geq 1 - 2 |\mathcal{H}| \, e^{-\frac{1}{2} \, m \, \varepsilon^2} = \underbrace{1 - \delta}_{\text{probably}}$$

• Given  $\varepsilon > 0$ , probability of failure  $\delta > 0$ , we can compute the sample complexity:

$$m_{ extsf{pac}}^{\mathcal{H}}(arepsilon,\delta) = rac{2}{arepsilon^2} \ln \left(rac{2\,|\mathcal{H}|}{\delta}
ight)$$

## **Probably Approximately Correct learning**

- Successful PAC learning algorithm: An algorithm can learn a hypothesis that is likely ("probably") to be approximately correct, given a sufficient number of training examples.
- **Definition:** Algorithm is a <u>successful PAC learner</u> for hypothesis class  $\mathcal{H}$  if there exists a function  $m_{\mathrm{pac}}^{\mathcal{H}} \colon (0,1) \times (0,1) \to \mathbb{N}$  such that: For every  $\varepsilon \in (0,1)$ ,  $\delta \in (0,1)$ , and every distribution p(x,y), when running the algorithm on  $m \geq m_{\mathrm{pac}}^{\mathcal{H}}(\varepsilon,\delta)$  examples  $T_m$  i.i.d. drawn from p(x,y), then the algorithm returns  $h_m = A(\mathcal{T}^m)$  such that

$$\mathbb{P}_{T_m \sim p^m} \Big( \underbrace{R(p, h_m) - R(p, h_{\mathcal{H}}) \leq \varepsilon}_{\text{approximately correct}} \Big) \geq \underbrace{1 - \delta}_{\text{probably}}$$

#### Key Concepts:

- Approximately correct: The learned predictor's risk is at most  $\varepsilon$  greater than the risk of the best possible predictor in the class  $\mathcal{H}$ .
- **Probably:** The probability of the algorithm failing to produce the approximately correct predictor is at most  $\delta$ .
- Sample complexity  $m_{\text{pac}}^{\mathcal{H}}(\varepsilon, \delta)$ : The minimum number of examples required to guaranteed the desired accuracy  $\varepsilon$  and the confidence  $1 \delta$ .
- **Distribution independence:** The guarantees hold for any data distribution p(x, y).



**Theorem.** Let  $\mathcal{H} = \{h^i : \mathcal{X} \to \mathcal{Y} \mid i \in \{1, \dots, H\}\}$  be a finite hypothesis class. Then the Empirical Risk Minimization (ERM) algorithm

$$h_m = \operatorname*{arg\,min}_{h \in \mathcal{H}} \left[ rac{1}{m} \sum_{i=1}^m \llbracket y_i 
eq h(x_i) 
rbracket 
ight]$$

is a successful PAC learner with sample complexity

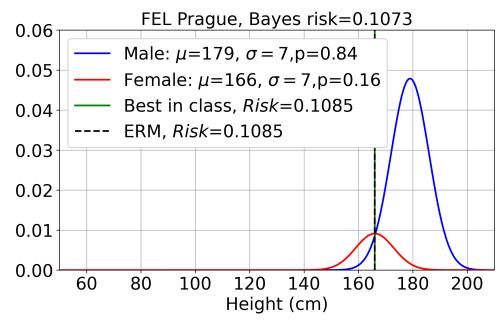
$$m_{\mathsf{PAC}}^{\mathcal{H}}(arepsilon,\delta) = rac{2}{arepsilon^2} \ln\!\left(\!rac{2\,|\mathcal{H}|}{\delta}\!
ight)\!.$$

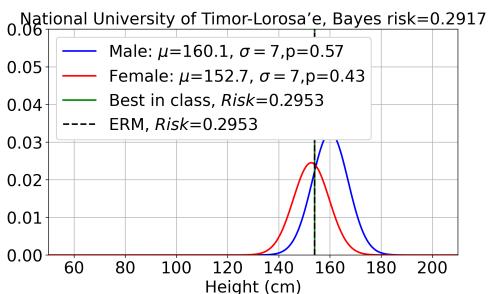
The theorem is a consequence of the bound we introduced earlier (however, have not yet proved):

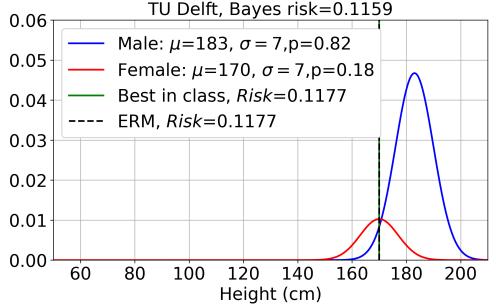
$$\mathbb{P}_{T_m \sim p^m} \left[ \underbrace{R(p, h_m) - R(p, h_{\mathcal{H}}) \leq \varepsilon}_{\text{approximately correct}} \right] \geq 1 - 2 |\mathcal{H}| \, e^{-\frac{1}{2} \, m \, \varepsilon^2} = \underbrace{1 - \delta}_{\text{probably}}$$

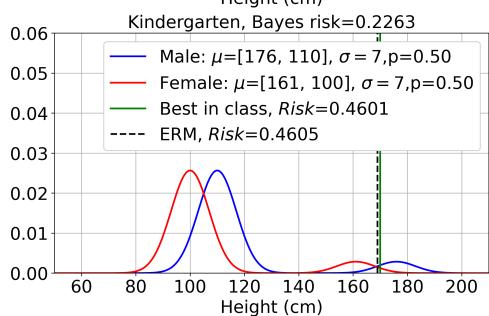


**Example:**  $\varepsilon = 0.05$ ,  $\delta = 0.1$ ,  $\mathcal{H} = \left\{ h(x; \theta) = \operatorname{sign}(x - \theta) \mid \theta \in \{140, 141, \dots, 200\} \right\}$ The sample complexity:  $m_{\mathsf{pac}}^{\mathcal{H}}(\varepsilon, \delta) = \frac{2}{\varepsilon^2} \ln \left( \frac{2 \mid \mathcal{H} \mid}{\delta} \right) = \frac{2}{0.05^2} \ln \left( \frac{2 \cdot 60}{0.1} \right) \approx 5,673$ 



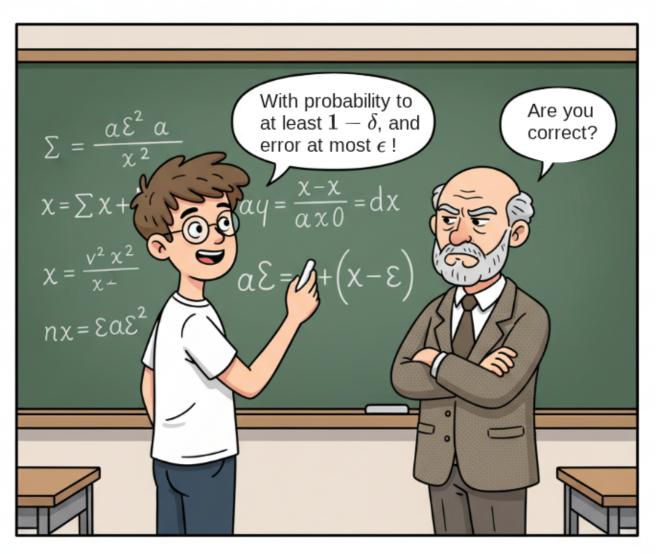






## **Probably Approximately Correct Joke**



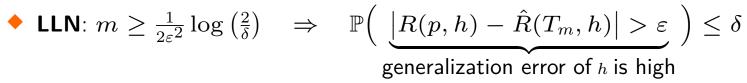


Joke: Chat GPT-5

Image: Nano Banana

BE4M33SSU exam edition.

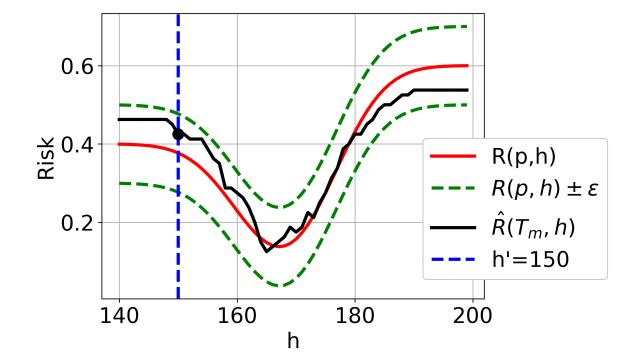
## **Uniform Law of Large Numbers**



LLN applies for any  $h\colon \mathcal{X} o \mathcal{Y}$  fixed prior to observating the data  $T_m$ 

$$\bullet \ \, \textbf{ULLN:} \ \, m \geq m_{\mathrm{ul}}^{\mathcal{H}}(\varepsilon,\delta) \quad \Rightarrow \quad \mathbb{P}\Big( \underbrace{\max_{h \in \mathcal{H}} \big| R(p,h) - \hat{R}(T_m,h) \big| > \varepsilon}_{\text{generalization error of some } h \in \mathcal{H} \text{ is high}} \Big) \leq \delta$$

ULLN applies only for some  $\mathcal{H}$ , e.g., when  $\mathcal{H}$  is finite  $m_{\mathsf{ul}}^{\mathcal{H}}(\varepsilon,\delta) = \frac{1}{2\varepsilon^2}\log\left(\frac{2|\mathcal{H}|}{\delta}\right)$ 



#### **Example:**

$$\mathcal{H} = \{h(x; \theta) = \text{sign}(x - \theta) \mid \theta \in \{140, 141, \dots, 200\}\}$$
 $\varepsilon = 0.1, \delta = 0.05, |\mathcal{H}| = 60$ 
 $m_{\text{ul}}^{\mathcal{H}}(0.1, 0.05) = 389.2$ 

#### 12/16

## **Uniform Law of Large Numbers**

• Assume a finite hypothesis class  $\mathcal{H} = \{h^i \colon \mathcal{X} \to \mathcal{Y} \mid i \in \{1, 2, \dots, H\}\}$ .

$$\mathbb{P}\Big(\underbrace{\max_{h\in\mathcal{H}}\left|R(p,h)-\hat{R}(T_m,h)\right|\geq\varepsilon}_{\text{generalization error of some }h\in\mathcal{H}\text{ is high}}\Big) \quad \stackrel{(1)}{=} \quad \mathbb{P}\left(\begin{array}{c} \left|R(p,h^1)-\hat{R}(T_m,h^1)\right|\geq\varepsilon & \text{or }\\ \left|R(p,h^2)-\hat{R}(T_m,h^2)\right|\geq\varepsilon & \text{or }\\ \vdots & \\ \left|R(p,h^H)-\hat{R}(T_m,h^H)\right|\geq\varepsilon \end{array}\right) \quad \stackrel{(2)}{\leq} \quad \sum_{h\in\mathcal{H}}\mathbb{P}\Big(\left|R(p,h)-\hat{R}(T_m,h)\right|\geq\varepsilon\Big) \quad \stackrel{(3)}{\leq} \quad 2\left|\mathcal{H}\right|e^{-2m\varepsilon^2}$$

- 1.  $a \ge \varepsilon$  or  $b \ge \varepsilon \iff \max\{a, b\} \ge \varepsilon$
- 2. Union bound:  $\mathbb{P}(A_1 \text{ or } A_2 \text{ or } \cdots \text{ or } A_n) \leq \sum_{i=1}^n \mathbb{P}(A_i)$
- 3. Hoeffding inequality:  $\mathbb{P}\Big( |R(h) R_{\mathcal{T}^m}(h)| \geq \varepsilon \Big) \leq 2 \, e^{-2 \, m \, \varepsilon^2}$
- Setting  $2 |\mathcal{H}| e^{-2m\varepsilon^2} = \delta$  and solving for m, we get

$$m_{\mathrm{ul}}^{\mathcal{H}}(\varepsilon,\delta) = \frac{1}{2\varepsilon^2} \log\left(\frac{2|\mathcal{H}|}{\delta}\right) \qquad \Rightarrow \qquad \mathbb{P}\Big(\max_{h\in\mathcal{H}} \left|R(h) - R_{\mathcal{T}^m}(h)\right| \geq \varepsilon\Big) \leq \delta$$

13/16

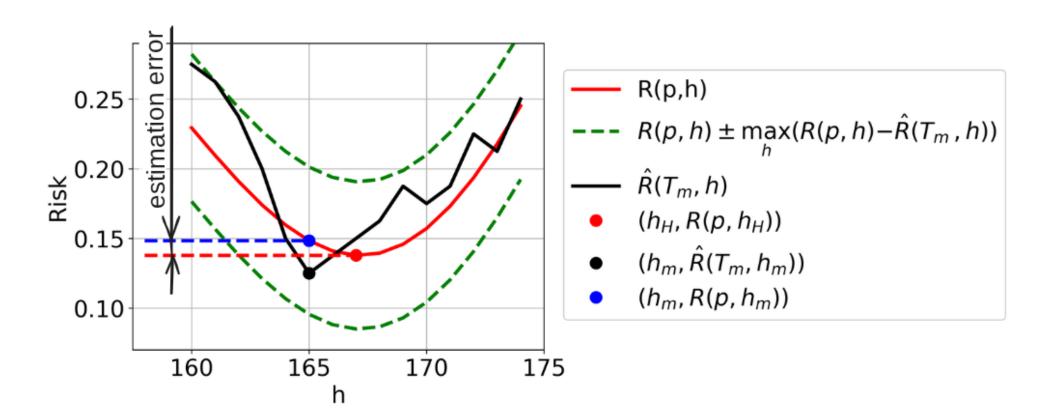
### **Bound on Estimation Error**

• Bound on estimation error of  $h_m = \arg\min_{h \in \mathcal{H}} \hat{R}(T_m, h)$  :

$$\underbrace{R(p,h_m) - R(p,h_{\mathcal{H}})}_{\textit{estimation error}} = \left(R(p,h_m) - \hat{R}(T_m,h_m)\right) + \left(\hat{R}(T_m,h_m) - R(p,h_{\mathcal{H}})\right)$$

$$\leq \left(R(p,h_m) - \hat{R}(T_m,h_m)\right) + \left(\hat{R}(T_m,h_{\mathcal{H}}) - R(h_{\mathcal{H}})\right)$$

$$\leq 2 \max_{h \in \mathcal{H}} \left|R(p,h) - \hat{R}(T_m,h)\right|$$
Maximal generalization error



## **ULLN** implies **ERM** is **PAC** learner

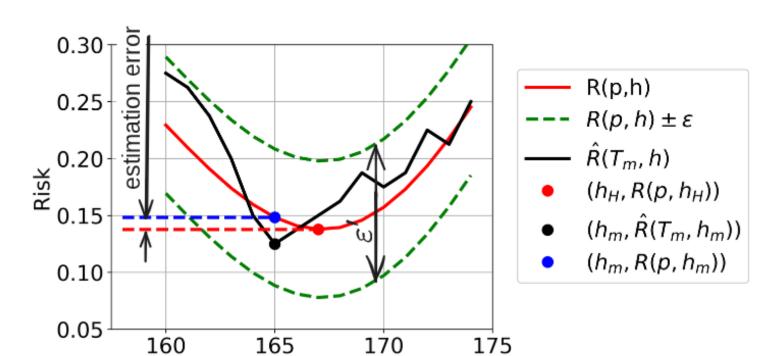


- ULLN:  $m \ge m_{\mathrm{ul}}^{\mathcal{H}}(\varepsilon, \delta) \quad \Rightarrow \quad \mathbb{P}\Big(\max_{h \in \mathcal{H}} \big| R(p, h) \hat{R}(T_m, h) \big| > \varepsilon \Big) \le \delta$
- Bound on estimation error of  $h_m = \arg\min_{h \in \mathcal{H}} \hat{R}(T_m, h)$ :

$$R(p, h_m) - R(p, h_{\mathcal{H}}) \le 2 \max_{h \in \mathcal{H}} |R(p, h) - \hat{R}(T_m, h)|$$

**♦** ULLN + Bound on estimation error = ERM is successfull PAC learner:

$$m \ge m_{\mathrm{pac}}^{\mathcal{H}}(\varepsilon', \delta) = m_{\mathrm{ul}}^{\mathcal{H}}(\varepsilon'/2, \delta) \quad \Rightarrow \quad \mathbb{P}\Big(R(p, h_m) - R(p, h_{\mathcal{H}}) \le \varepsilon'\Big) \ge 1 - \delta$$



h

#### **Example:**

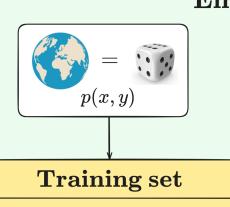
For finite  $\mathcal{H}$ , we have:

$$m_{\mathrm{ul}}^{\mathcal{H}}(arepsilon,\delta) = rac{1}{2arepsilon^2}\log\left(rac{2|\mathcal{H}|}{\delta}
ight)$$

Hence, the sample complexity is:

$$m_{ extsf{pac}}^{\mathcal{H}}(arepsilon',\delta) = rac{2}{arepsilon'^2}\log\left(rac{2|\mathcal{H}|}{\delta}
ight)$$





#### Hypothesis class

$$\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}} = \{h : \mathcal{X} 
ightarrow \mathcal{Y}\}$$

#### Loss function

$$\ell \colon \mathcal{Y} imes \mathcal{Y} o \mathbb{R}$$

i.i.d. data

$$igg| T_m = ((x_1,y_1),\ldots,(x_m,y_m)) \sim p^m$$

#### ERM Learning Algorithm

Learning algorithm  $A: (\mathcal{X} \times \mathcal{Y})^m \to \mathcal{H}$  finding a predictor in  $\mathcal{H}$  which minimizes the empirical error:

$$\hat{R}(T_m,h) = rac{1}{m} \sum_{i=1}^m \ell(y_i,h(x_i; heta))$$



#### **Predictor**

$$\hat{y} = h_m(x)$$

#### Error decomposition

Predictor Error = Estimation Error + Approximation Error + Bayes Error "Too complex" hypothesis space

E.g. Memorizer ERM is not PAC learner

#### Uniform Law of Large Numbers

ULLN holds for  $\mathcal{H} \iff$  ERM is PAC learner.

#### **PAC** learning

Successful PAC learner: finds approximately correct predictor with high probability.

$$m \geq m_{\mathrm{pac}}^{\mathcal{H}}(arepsilon, \delta)$$
:

 $\mathbb{P}[ ext{estimation error} \leq arepsilon] \geq 1 - \delta$ 

## Finite hypothesis space

$$\mathcal{H} = \{h_1, h_2, \ldots, h_{\mathcal{H}}\}$$

**ERM** is PAC learner

$$m_{ ext{pac}}^{\mathcal{H}}(arepsilon,\delta) = rac{2}{arepsilon^2} ext{log}\left(rac{2|\mathcal{H}|}{\delta}
ight)$$

#### **VC** dimension

 $\text{VCdim: } \{-1, +1\}^{\mathcal{X}} \to \mathbb{N}$ 

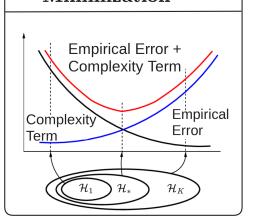
**Fundamental Theorem:** 

 $VCdim(\mathcal{H}) < \infty$ 

 $\iff$  ULLN holds for  ${\cal H}$ 

← ERM is PAC learner

#### Structured Risk Minimization



## **Summary of Key Concepts**

#### Error Decomposition

• Error of the learned predictor = Estimation Error + Approximation Error + Bayes Risk

#### Probably Approximately Correct (PAC) Learning

- A successful PAC learner, with high probability, finds a close approximation of the best predictor in the class, given enough examples.
- Sample complexity: number of examples needed for PAC guarantees.

#### Empirical Risk Minimization:

• ERM over a finite hypothesis space is a successful PAC learner.

#### Uniform Law of Large Numbers

• Guarantees uniform convergence of empirical risk to the expected risk over the hypothesis space.