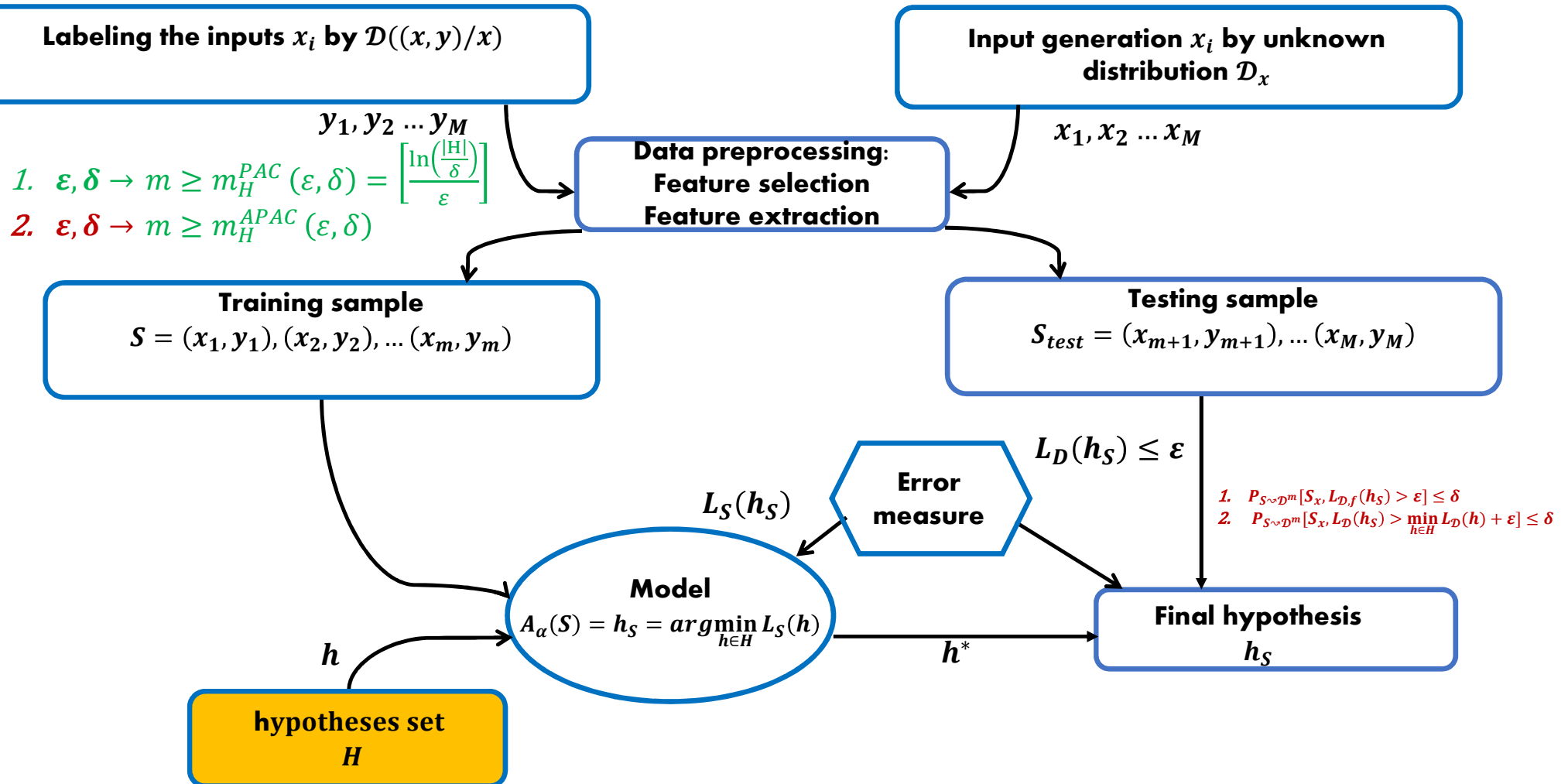


# Part 1: Machine learning theory

1. Learning framework
2. **Uniform convergence:**
  1.  $\epsilon$ -representative sample.
  2. Uniform convergence.
3. Learnability of infinite size hypotheses classes
4. Tradeoff Bias/Variance
5. Non-Uniform learning.

# Reminder



## Reminder

- **PAC** if hypotheses of realizability: target  $f$  exist
- **APAC** if there isn't hypotheses of realizability  $f$  with probability
- If  $|H| < \infty$  and hypotheses of realizability holds then we have PAC  $\forall m \geq m_H^{PAC}(\epsilon, \delta) = \left\lceil \frac{\ln(|H|)}{\epsilon} \right\rceil$

Learning PAC :  $m_H^{PAC}(\epsilon, \delta)$

- $\forall \epsilon, \delta \in [0,1]^2$ , and  $\forall \mathcal{D}$  over  $Z$ ,  $\exists m_H(\epsilon, \delta)$  such that  $\forall m > m_H^{PAC}(\epsilon, \delta)$  we have  $P_{S \sim \mathcal{D}^m}[S_x, L_{\mathcal{D},f}(h_S) > \epsilon] \leq \delta$

Learning APAC:  $m_H^{APAC}(\epsilon, \delta)$

- $\forall \epsilon, \delta \in [0,1]^2$  and  $\forall \mathcal{D}$  over  $Z$ ,  $\exists m_H(\epsilon, \delta)$  such that  $\forall m > m_H^{APAC}(\epsilon, \delta)$  we have  $P_{S \sim \mathcal{D}^m}[S_x, L_{\mathcal{D}}(h_S) > \min_{h \in H} L_{\mathcal{D}}(h) + \epsilon] \leq \delta$

## Motivation

Our aim in this chapter is to prove this proposition (there isn't a realizability hypotheses ):

**if  $H$  is a finite class of hypotheses, then  $H$  follows an Agnostic PAC learning.**

**Tool:**

Uniform convergence.

$|H| < \infty \implies \text{Learning Uniforme convergence} \implies \text{Agnostic PAC learning}$

**Learning Uniforme convergence:  $m_H^{CU}(\epsilon, \delta)$ :  $S$  is  $\epsilon$ -representative**

- $\forall (\epsilon, \delta) \in [0,1]^2$  and  $\forall \mathcal{D}$  over  $Z$ ,  $\exists m_H^{CU}(\epsilon, \delta)$  such that  $\forall m > m_H^{CU}(\epsilon, \delta)$

$$P[S_x: |L_S(h_S) - L_D(h_S)| > \epsilon] \leq \delta \Leftrightarrow P[S_x: |L_S(h_S) - L_D(h_S)| \leq \epsilon] \geq 1 - \delta$$

- $|L_S(h_S) - L_D(h_S)| \approx 0$  and  $L_D(h_S) \approx 0$

## 2.1. $\varepsilon$ -representative sample

### Definition:

The sample  $S \subset Z$  is  $\varepsilon$ -representative with respect to  $(Z, H, l, \mathcal{D})$  if :

$$\forall h \in H \quad |L_S(h) - L_D(h)| \leq \varepsilon$$

### Notice:

If  $S$  is  $\varepsilon$ -representative, so  $ERM_H$  is a good learning strategy.

- $|L_S(h) - L_D(h)| \approx 0$
- $L_D(h) \approx 0$

## 2.1. $\varepsilon$ -representative sample

### Lemma:

If  $S$  is  $\varepsilon$ -representative with respect to  $(Z, H, l, \mathcal{D})$ , then:

$$L_D(h_s) \leq \min_{h \in H} (L_D(h)) + 2\varepsilon$$

Such that:

$$ERM_H(S) = A_\alpha(S) = h_S \in \operatorname{argmin}_{h \in H} \{L_S(h)\}$$

## 2.1. $\varepsilon$ -representative sample

**proof:**

Let  $S$  be  $\varepsilon$ -representative, then:

$\forall h \in H$ , we have:

$$|L_S(h) - L_D(h)| \leq \varepsilon \Rightarrow L_S(h) \leq L_D(h) + \varepsilon$$

We know that  $h_S$  is the output of  $ERM_H(S)$ , so:

$$h_S \in \operatorname{argmin}_{h \in H} \{L_S(h)\}$$

So,  $\forall h \in H$ , we have:

$$L_S(h_S) \leq L_S(h)$$

Since  $S$  is  $\varepsilon$ -representative, then for  $h = h_S$ , We also have:

$$L_D(h_S) \leq L_S(h_S) + \varepsilon$$

## 2.1. $\varepsilon$ -representative sample

proof:

$$L_D(h_S) \leq L_S(h_S) + \varepsilon \leq L_S(h) + \varepsilon \leq L_D(h) + 2\varepsilon$$

So,  $\forall h \in H$  :

$$L_D(h_S) \leq L_D(h) + 2\varepsilon$$

Then :

$$L_D(h_S) \leq \min_{h \in H} (L_D(h)) + 2\varepsilon$$



## 2.2. Uniform convergence

### Definition:

We say that  $H$  has the uniform convergence property with respect to  $(Z, l)$ , if there exist:

- a function  $m_H^{CU}: [0,1]^2 \rightarrow \mathbb{N}$ , such that:  $\forall (\varepsilon, \delta) \in [0,1]^2$  and  $\forall \mathcal{D}$  over  $Z$ .
- $S$  is a sample of size  $m \geq m_H^{CU}(\varepsilon, \delta)$ , whose points are drawn **(i.i.d.)** by  $\mathcal{D}$ , such that with probability of at least  $(1 - \delta)$ ,  $S$  is  $\varepsilon$ -representative (avec probability  $1 - \delta$ ):

$$\text{eq 1: } P[S_x: |L_S(h_S) - L_D(h_S)| \leq \varepsilon] \geq 1 - \delta \Leftrightarrow P[S_x: |L_S(h_S) - L_D(h_S)| > \varepsilon] \leq \delta$$

$P[S_x: |L_S(h_S) - L_D(h_S)| > \varepsilon] \leq \delta$  bad hypothesis (there isn't a  $\varepsilon$ -representative )

Such that  $ERM_H(S) = A_\alpha(S) = h_S \in \underset{h \in H}{\operatorname{argmin}} \{L_S(h)\}$

### Notice:

If  $H$  has the uniform convergence property  $\Rightarrow H$  is called « **Glivenko-Cantelli class** ».

H: hypotheses set

**PAC:**  $m_H^{PAC}(\epsilon, \delta), \epsilon_{PAC}$   
 $P_{S \sim \mathcal{D}^m}[S_x, L_{\mathcal{D},f}(h) > \epsilon] \leq \delta$

**APAC**  $m_H^{APAC}(\epsilon, \delta)$   
 $P_{S \sim \mathcal{D}^m}[L_{\mathcal{D}}(h_S) > \min_{h \in H} L_{\mathcal{D}}(h) + \epsilon_{APAC}] \leq \delta$

**CU:**  $m_H^{CU}(\epsilon_{CU}, \delta)$   
 $P_{S \sim \mathcal{D}^m}[|L_S(h) - L_D(h)| > \epsilon] \leq \delta$   
***S* is  $\epsilon$ -representative**

## 2.2. Uniform convergence

**Definition:** Sample complexity  $m_H^{CU}(\varepsilon, \delta) \Rightarrow m_H^{APAC}(2\varepsilon, \delta) \Rightarrow m_H^{PAC}(\varepsilon, \delta)$

The function  $m_H^{CU}(\varepsilon, \delta): [0,1]^2 \rightarrow \mathbb{N}$  enables to determine the minimal number of data for which  $H$  follows a uniform convergence with accuracy  $\varepsilon$  and confidence  $\delta$ .

### Theorem1: (Lemma 1 and definition Of CU)

If  $H$  follows a uniform convergence with complexity  $m_H^{CU}(\varepsilon, \delta)$ , then  $H$  follows Agnostic PAC learning with complexity  $m_H^{APAC}(\varepsilon, \delta)$  such that:

$$m_H^{APAC}(\varepsilon, \delta) \leq m_H^{CU}(\varepsilon/2, \delta) \Rightarrow m_H^{APAC}(\varepsilon, \delta) \approx m_H^{CU}(\varepsilon/2, \delta)$$

Moreover,  $ERM_H$  succeeds in the agnostic PAC learning of  $H$ .

## 2.2. Uniform convergence

### General Objective

#### Objective of proof:

Prove that if  $H$  is a finite class of hypotheses  $\Rightarrow H$  is agnostic PAC learnable.

But, we know that:

### Remarks

If  $S$  is  $\varepsilon$ -representative  $\Rightarrow H$  has the uniform convergence property  $\Rightarrow H$  is agnostic PAC learnable.(theorem 1 and Lemma)

#### Objective: reformulation

We should simply prove that, if  $|H| < \infty$  and we have sufficient amount of data,  $S$  is  $\varepsilon$ -representative (avec probability  $1 - \delta$ ).

Then  $H$  is agnostic PAC learnable(Remarks)

## 2.2. Uniform convergence

### Proof strategy of the general objective

#### Step 1:

Consider that  $H$  owns one hypothesis  $h$ , and let's prove that  $S$  is  $\varepsilon$ - representative:

$$L_{\mathcal{D}}(h) \approx L_S(h)$$

#### Step 2:

Let's suppose that  $H$  owns many hypotheses, and let's prove using Boole's inequality that  $S$  is  $\varepsilon$ - representative:

$$\forall h \in H, \quad L_{\mathcal{D}}(h) \approx L_S(h)$$

#### Step 3:

Let's determine the necessary number of data so that  $S$  can be  $\varepsilon$ - representative.

## 2.2. Uniform convergence

### Definition: Law of large scale numbers

Let  $(\theta_1, \theta_2, \dots, \theta_m)$  be a random variables (**i. i. d.**), such that  $\mu$  is their real mean and  $\frac{1}{m} \sum_{i=1}^m \theta_i$  their empirical mean.  $\left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| \rightarrow 0$  when  $m \rightarrow \infty$

So, if  $m \rightarrow \infty$ , the empirical mean converges converge to the real mean, with probability equal to 1.

The measures of concentration inequalities are statistical tools that allow to quantify the deviation between the empirical mean and the real mean when  $m$  is finite.

Among these inequalities, there exists “**Hoeffding’s inequality**”.

## 2.2. Uniform convergence

### Definition: Hoeffding's Inequality

Let's suppose that  $(\theta_1, \theta_2, \dots, \theta_m)$  are random variables (*i.i.d.*), having the real mean  $\mu$ , such that these variables have values in  $[a, b]$ . So:

$$P \left[ \left| \mu - \frac{1}{m} \sum_{i=1}^m \theta_i \right| > \epsilon \right] \leq 2e^{\frac{-2m\epsilon^2}{(b-a)^2}}$$

Real mean  
( $L_{\mathcal{D}}$  general)

$L_S$  Empirical mean

### Notice:

This probability decreases exponentially if the size  $m$  of the sample increases.

$$P \left[ \left| (\mu = L_D(h)) - \left( \frac{1}{m} \sum_{i=1}^m \theta_i = L_S(h) \right) \right| > \epsilon \right] \leq \delta$$

## 2.2. Uniform convergence

### Proof - Step 1

Let's suppose that  $H = \{h\}$ , and let's prove that  $S$  is  $\varepsilon$ - representative:

$$|L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon$$

This implies to prove that :

$$\mathbb{P}_{S \sim \mathcal{D}^m}[|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] \text{ is small}$$

To bound this inequality we will use the Hoeffding's Inequality.

We have:

$$\mathbb{P}_{S \sim \mathcal{D}^m}[|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] = P \left[ \left| \mathbb{E}_{z \sim \mathcal{D}} [l(h, z)] - \frac{1}{m} \sum_{i=1}^m l(h, z_i) \right| > \varepsilon \right]$$



## 2.2. Uniform convergence

### Proof - Step 1

In that case we have:

$$\theta_i = l(h, z_i) \in [0,1] \quad , \quad \mu = \mathbb{E}_{z \sim \mathcal{D}}[l(h, z)] \quad , \quad \frac{1}{m} \sum_{i=1}^m \theta_i = \frac{1}{m} \sum_{i=1}^m l(h, z_i)$$

Then:

$$P \left[ \left| \mathbb{E}_{z \sim \mathcal{D}}[l(h, z)] - \frac{1}{m} \sum_{i=1}^m l(h, z_i) \right| > \epsilon \right] \leq 2e^{-2m\epsilon^2}$$

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} [|L_{\mathcal{D}}(\mathbf{h}) - L_{\mathcal{S}}(\mathbf{h})| > \epsilon] \leq 2e^{-2m\epsilon^2} \quad \text{Eq.1}$$

So:

$\exists h \in H, \quad |L_{\mathcal{D}}(h) - L_{\mathcal{S}}(h)| \leq \epsilon$  if  $m$  is sufficiently large.

## 2.2. Uniform convergence

### Proof - Step 2

Now, let's generalize **Eq.1** for all hypotheses  $h \in H$ .

Let's suppose that  $H$  owns several hypotheses, and let's prove that using the Boole's inequality that  $S$  is  $\varepsilon$ -representative:

$$\forall h \in H, \quad L_{\mathcal{D}}(h) \approx L_S(h)$$

We have:

$$(\forall h \in H, L_{\mathcal{D}}(h) \approx L_S(h)) \Leftrightarrow (\forall h \in H, \text{the probability of failure of } h \text{ is small})$$

We have that  $h$  fails if  $S$  is not  $\varepsilon$ -representative.

This implies to prove that  $\forall h \in H$ :

$$P[S \text{ is not } \varepsilon\text{-representative with respect to } H] \text{ is small}$$

## 2.2. Uniform convergence

### Proof - Step 2

$$\begin{aligned} P[S \text{ is not } \varepsilon\text{-representative with respect to } H] \\ = P[\exists h \in H, |L_D(h) - L_S(h)| > \varepsilon] \end{aligned}$$

We have:

$$P[\exists h \in H, |L_D(h) - L_S(h)| > \varepsilon] \leq P\left[\bigcup_{h \in H} \{h, |L_D(h) - L_S(h)| > \varepsilon\}\right]$$

According to **Boole's inequality**:

$$P\left[\bigcup_{h \in H} \{h, |L_D(h) - L_S(h)| > \varepsilon\}\right] \leq \sum_{h \in H} P_{S \sim \mathcal{D}^m}[|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon]$$

According to Hoeffding inequality, we have:

$$P_{S \sim \mathcal{D}^m}[|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] \leq 2e^{-2m\varepsilon^2}$$

## 2.2. Uniform convergence

### Proof - Step 2

So:

$$\sum_{h \in H} P_{S \sim \mathcal{D}^m} [|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] \leq \sum_{h \in H} 2e^{-2m\varepsilon^2}$$

Then we will have that:

$$P[S \text{ is not } \varepsilon\text{-representative with respect to } H] \leq |H|2e^{-2m\varepsilon^2}$$

So:

$$\forall h \in H \quad P_{S \sim \mathcal{D}^m} [|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] \leq |H|2e^{-2m\varepsilon^2} \quad \text{Eq.2}$$

Finally:

$\forall h \in H, L_{\mathcal{D}}(h) \approx L_S(h)$  if  $m$  is sufficiently big.

## 2.2. Uniform convergence

### Proof - Step 3

Let's determine the necessary number of data so that  $S$  can be  $\varepsilon$ -representative.

We know that  $m_H^{CU}(\varepsilon, \delta)$  is the minimal number of data so that  $S$  can be  $\varepsilon$ -representative with probability  $\geq 1 - \delta$ .

So, we want that  $P[S \text{ is not } \varepsilon\text{-representative with respect to } H] \leq \delta$

So:

$$P_{S \sim \mathcal{D}^m}[|L_{\mathcal{D}}(h) - L_S(h)| > \varepsilon] \leq |H|2e^{-2m\varepsilon^2} \leq \delta$$

Hereby, the necessary number of data is:

$$m \geq \frac{\ln\left(\frac{2|H|}{\delta}\right)}{2\varepsilon^2}$$

## 2.2. Uniform convergence

### Theorem2:

Let  $H$  be a finite class of hypotheses, let  $Z$  be a set of data and  $l: H \times Z \rightarrow [0,1]$  the cost function.

So,  $H$  follows a uniform convergence learning with sample complexity:

$$m_H^{CU}(\epsilon, \delta) \leq \left\lceil \frac{\ln\left(\frac{2|H|}{\delta}\right)}{2\epsilon^2} \right\rceil \Rightarrow m_H^{CU}(\epsilon, \delta) \approx \left\lceil \frac{\ln\left(\frac{2|H|}{\delta}\right)}{2\epsilon^2} \right\rceil \text{ (by proof above )}$$

Moreover,  $H$  follows agnostic PAC learning with  $ERM_H$  algorithm, with sample complexity:

$$m_H^{APAC}(\epsilon, \delta) \leq m_H^{CU}\left(\frac{\epsilon}{2}, \delta\right) \leq \left\lceil \frac{2 \ln\left(\frac{2|H|}{\delta}\right)}{\epsilon^2} \right\rceil \text{ (by th1)}$$

$$\bullet \Rightarrow m_H^{APAC}(\epsilon, \delta) \approx m_H^{CU}\left(\frac{\epsilon}{2}, \delta\right) \approx \left\lceil \frac{2 \ln\left(\frac{2|H|}{\delta}\right)}{\epsilon^2} \right\rceil$$

# Supervised Learning Passive Offline Algorithm (SLPOA)

