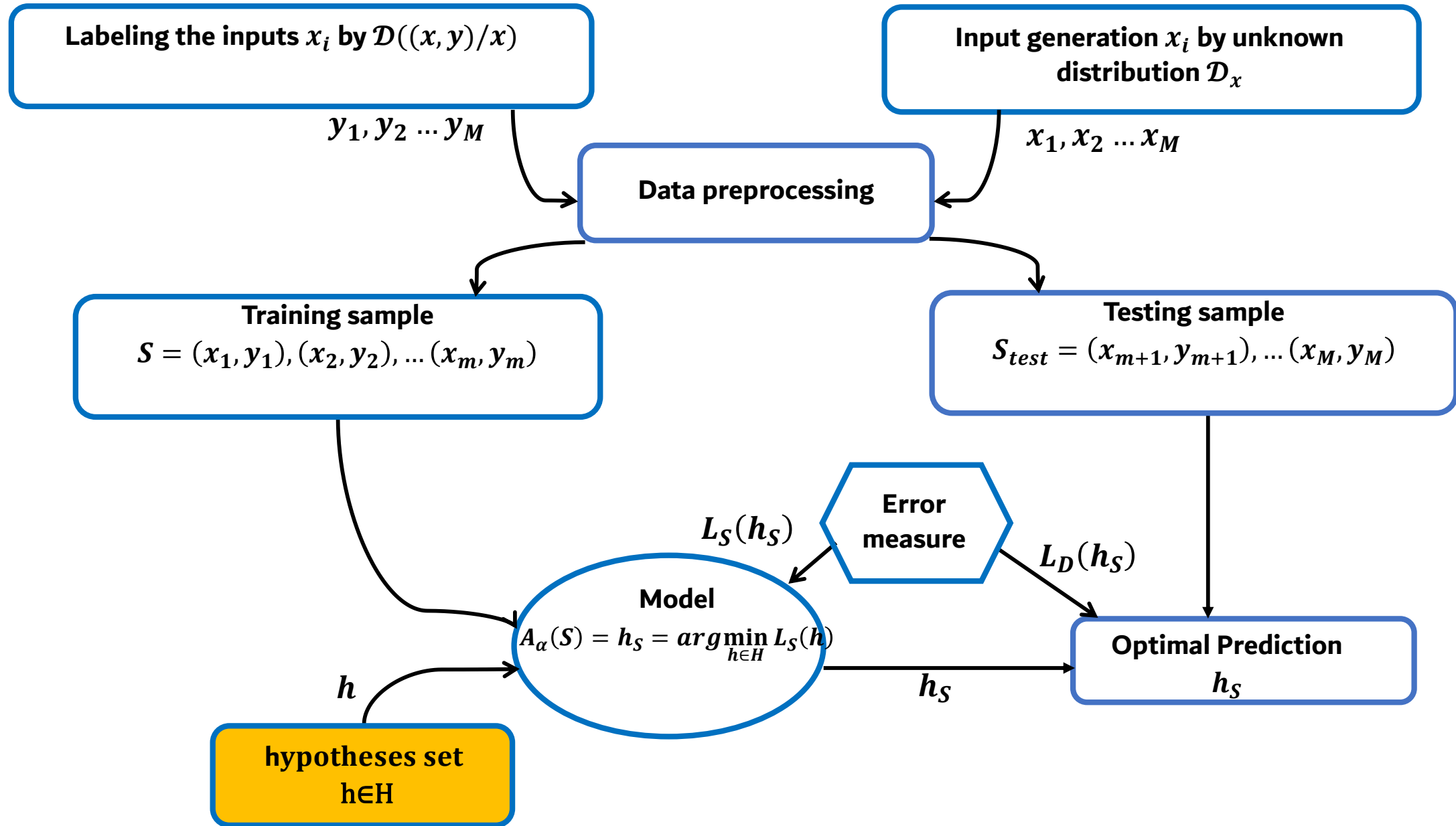


# Part 1: Machine learning theory

1. Learning framework
2. Uniform convergence
3. Learnability of infinite size hypotheses set
  1. No-Free-Lunch theorem
  2. Infinite hypothesis class: Exemple
  3. VC dimension
  4. Covering number
4. Tradeoff Bias/Variance
5. Non-Uniform learning

# Supervised Learning Passive Offline Algorithm (SLPOA)



# Recall (classification)

## Definition: shuttering

Let  $\mathbf{H}$  be a set of functions from  $X$  to  $\{0,1\}$  and  $A \subseteq X$  a finite set.

We say that  $\mathbf{H}$  shutteres  $A$  if the restriction of  $H$  over  $A$  is of finite cardinality:

$$|\mathbf{H}_A| = 2^{|A|}$$

Such that:

$$\mathbf{H}_A = \{h(a_1), \dots, h(a_{|A|}): h \in \mathbf{H}\}$$

# Recall (classification)

## Definition: VC Dimension

The VC dimension is a property of  $H$  which measures the maximum size of a set  $A$  to be shattered by  $H$ :

$$d_{VC}(H) = \begin{cases} \max\{|A|, A \text{ is shattered by } H\} \\ +\infty \text{ there is no maximum for } A \end{cases}$$

## Definition: Growth function

Let  $H$  be a class of hypothesis, the **growth function** of  $H$  is  $\Pi_H: \mathbb{N} \rightarrow \mathbb{N}$ , such that:

$$\Pi_H(m) = \max_{\substack{A \subset X \\ |A|=m}} |H_A| \quad H_A \text{ is the restriction of } H \text{ on } A.$$

# Recall (classification)

## Theorem:

Let  $H$  be a class of hypotheses such that  $d_{VC}(H) = +\infty$ . So  $H$  is not PAC.

## Notice:

- $\forall H$  and  $\forall m, \Pi_H(m) \leq 2^m$
- If  $H$  shatters the class of size  $m$ , So:  $\Pi_H(m) = 2^m$
- **If  $d_{VC}(H) = \max\{|A|, A \text{ is shattered by } H\} < m$ , So:  $\Pi_H(m) < 2^m$**

## Lemma:

$\forall H$  and  $\forall A \subseteq X: |H_A| \leq |\{B \subseteq A : H \text{ shatters } B\}|$

## Lemma: Sauer

Let  $H$  be a class of hypotheses such that:  $d_{VC}(H)(\approx) \leq d < +\infty$

Then:  $\forall m, \Pi_H(m)(\approx) \leq \sum_{i=0}^d C_m^i \Rightarrow \log \Pi_H(m) \leq \log \sum_{i=0}^d C_m^i$

In particular, if  $m > d + 1$ , so:  $\Pi_H(m)(\approx) \leq \left(\frac{me^1}{d}\right)^d \Rightarrow \log \Pi_H(m) \leq \log \left(\frac{me^1}{d}\right)^d$

$$\Rightarrow 4 + \sqrt{\log(\Pi_H(m))} \leq 4 + \sqrt{\log \left(\frac{me^1}{d}\right)^d}$$

# Recall(classification)

## Theorem: Generalization bound of VC

Let  $H$  be a class of hypotheses and  $\Pi_H$  is its growth function. So, for any  $D$  and for any  $\delta \in [0,1]$ :

$$P_{S \sim D^m}(|L_D(h) - L_S(h)| \leq \varepsilon) \geq 1 - \delta$$

Such that:  $\varepsilon = \frac{4 + \sqrt{\log(\Pi_H(2m))}}{\delta \sqrt{2m}}$

## Lemma: Sauer

Let  $H$  be a class of hypotheses such that:  $d_{VC}(H)(\approx) \leq d < +\infty$

Then:  $\forall m, \quad \Pi_H(m)(\approx) \leq \sum_{i=0}^d C_m^i \Rightarrow \log \Pi_H(m) \leq \log \sum_{i=0}^d C_m^i$

In particular, if  $m > d + 1$ , so:  $\Pi_H(m)(\approx) \leq \left(\frac{me^1}{d}\right)^d \Rightarrow \log \Pi_H(m) \leq \log \left(\frac{me^1}{d}\right)^d$

$$\Rightarrow 4 + \sqrt{\log(\Pi_H(m))} \leq 4 + \sqrt{\log \left(\frac{me^1}{d}\right)^d}$$

# Recall(classification)

## Theorem:

Let  $H$  be a class of hypotheses in  $X \times \{0,1\}$ .

Let  $l$  be the classification loss function.

We have equivalence between:

1.  $H$  follows a uniform convergence.
2.  $H$  is agnostic PAC learnable by ERM.
3.  $H$  is agnostic PAC learnable.
4.  $H$  is PAC learnable.
5.  $H$  is PAC learnable by ERM.
6.  $d_{VC}(H)$  is finite.

## Notice:

The VC dimension is a tool characterizing the PAC learning.

# Covering number

1. Background
2. Covering numbers in a general metric space
3. Covering numbers in Euclidean space
4. Uniform convergence in a Real-valued Function class  $H$



# 1. Background

## Definition: Metric space

$(M, d)$  is called a metric space that consists of a set  $M$  together with a metric  $d: M \times M \rightarrow [0, \infty)$  that satisfies the following for all  $x, y, z \in M$ :

- $d(x, y) = 0 \implies x = y.$
- $d(x, y) = d(y, x).$
- $d(x, z) \leq d(x, y) + d(y, z).$

## Definition: Open $d$ -ball

An open  $d$ -ball centered at  $x \in M$  is defined as:

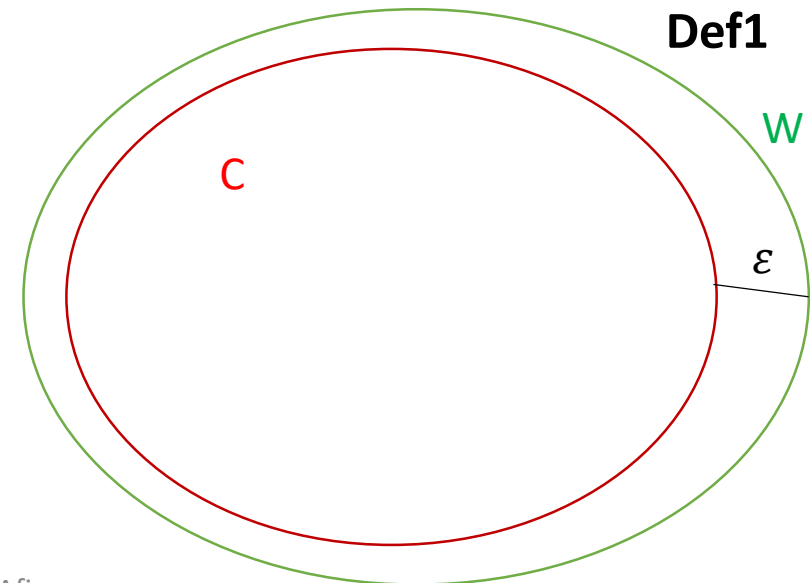
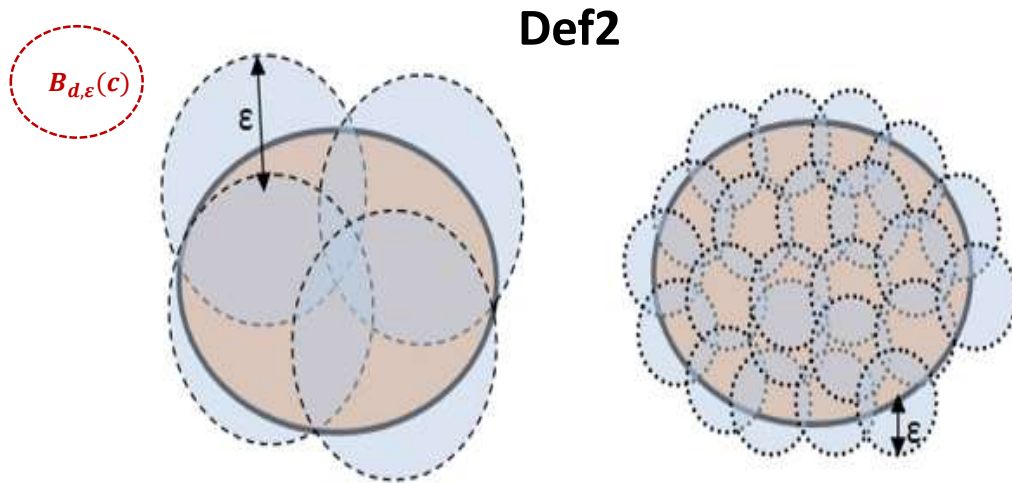
$$B_{d,\varepsilon}(x) = \{y \in M \mid d(x, y) < \varepsilon\}$$

## 2. Covering numbers in a general metric space

### Definition: $\varepsilon$ -cover

Let  $(M, d)$  be a metric space. Let  $W \subseteq M$  and let  $\varepsilon > 0$

- **Def1:** A set  $C \subseteq W$  is said to be  $\varepsilon$ -cover of  $W$  with respect to  $d$  if  $(\forall w \in W)(\exists c \in C)$  such that:  $d(w, c) < \varepsilon$
- **Def2:**  $C \subseteq W$  is an  $\varepsilon$ -cover of  $W$  with respect to  $d$  if the union of (open)  $d$ -balls of radius  $\varepsilon$  centered at points in  $C$  contains  $W$ :  $\bigcup_{c \in C} B_{d,\varepsilon}(c) \supseteq W$



## 2. Covering numbers in a general metric space

### Definition: $\varepsilon$ -covering number

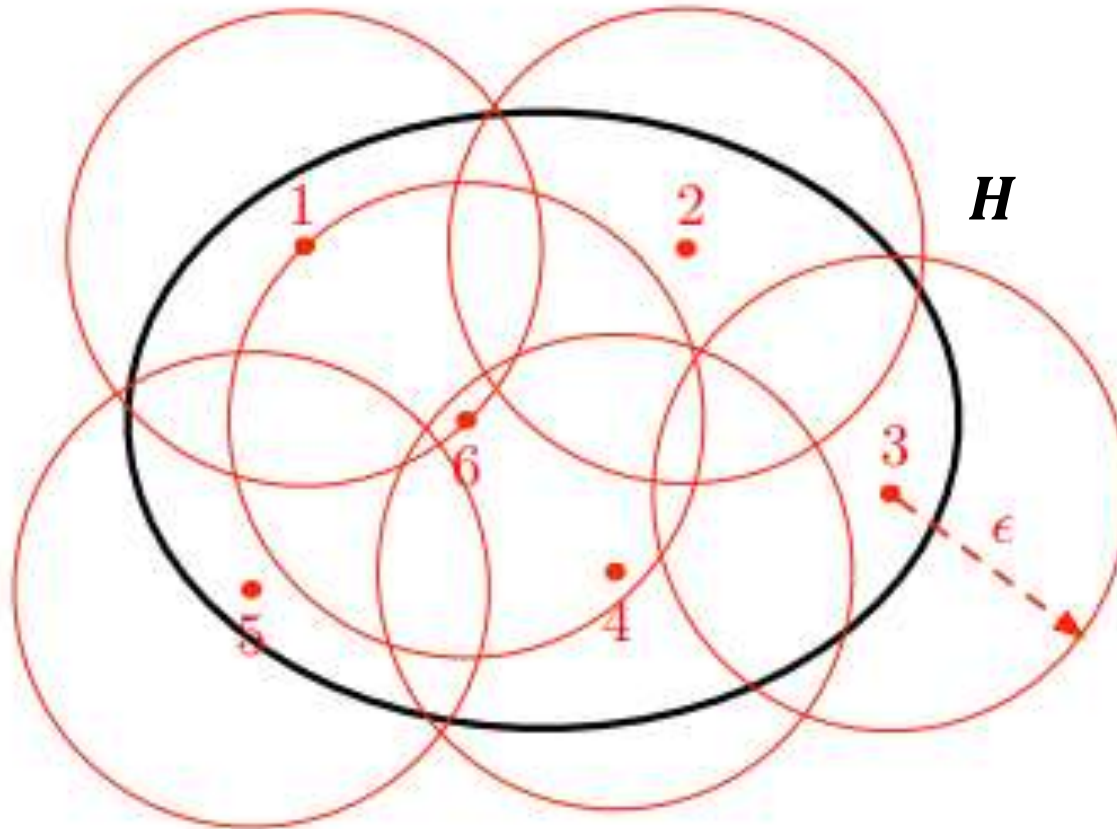
The  $\varepsilon$ -covering number  $\mathcal{N}(\varepsilon, W, d)$  of  $W$  with respect to  $d$  is defined as the cardinality of the smallest  $\varepsilon$ -cover of  $W$  if  $W$  has a finite  $\varepsilon$ -cover with respect to  $d$ . Otherwise, if  $W$  does not have a finite  $\varepsilon$ -cover with respect to  $d$ ,  $\varepsilon$ -covering number is equal to infinity.

$$\mathcal{N}(\varepsilon, W, d) = \begin{cases} \min\{|C|, C \text{ is an } \varepsilon - \text{cover of } W \text{ with respect to } d\} \\ \infty & \text{if } W \text{ does not have a finite } \varepsilon - \text{cover} \end{cases}$$

## 2. Covering numbers in a general metric space

### Example:

For instance, for the  $H$  shown in the figure the set of points  $\{1, 2, 3, 4, 5, 6\}$  is a covering. However, the covering number is 5 as point 6 can be removed from the set  $C$  and the resulting points are still a covering.  $C = \{1, 2, 3, 4, 5\}$



### 3. Covering numbers in Euclidean space

Consider now  $M = \mathbb{R}^n$ . We can define a number of different metrics on  $\mathbb{R}^n$ , including in particular the following:

$$d_1(x, y) = \frac{1}{n} \sum_{i=1}^n |x_i - y_i|$$

$$d_2(x, y) = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2}$$

$$d_\infty(x, y) = \max_i |x_i - y_i|$$

### 3. Covering numbers in Euclidean space

Accordingly, for any  $W \subseteq \mathbb{R}^n$ , we can define the corresponding covering numbers  $\mathcal{N}(\varepsilon, W, d)$  for  $p = 1, 2, \infty$ .

It is easy to see that:

$$d_1(x, y) \leq d_2(x, y) \leq d_\infty(x, y)$$

Therefore, the corresponding covering numbers satisfy the relation:

$$\mathcal{N}(\varepsilon, W, d_1) \leq \mathcal{N}(\varepsilon, W, d_2) \leq \mathcal{N}(\varepsilon, W, d_\infty)$$

## 4. Uniform covering numbers for a real-valued function class $H$

### Definition: uniform covering number

Let  $H$  be a class of real-valued functions on  $X$ :

$$H = \{h \mid h: X \rightarrow \mathbb{R}\}$$

And let  $S = \{x_1, \dots, x_m\} \subset X$ . Then the  $H(S) = H_S \subseteq \mathbb{R}$ .

For any  $\varepsilon > 0$  and  $m \in \mathbb{N}$ , the uniform  $d_p$  covering numbers of  $H$  for  $p = 1, 2, \infty$  are defined as:

$$\mathcal{N}_p(\varepsilon, H, m) = \begin{cases} \max_{S \subset X: |S|=m} \mathcal{N}(\varepsilon, H_S, d_p) & \text{if } \mathcal{N}(\varepsilon, H_S, d_p) \text{ is finite for all } S \subset X \\ \infty & \text{otherwise} \end{cases}$$

**Notice:** The number of “uniform” refers to the maximum over all  $S \subset X$ . It has no relationship with uniform convergence.

## 4. Uniform covering numbers for a real-valued function class $H$

### Background:

- $S = \{(x, y)\} \subseteq X \times Y$ : **is data set with**  $Y \subseteq \mathbb{R}$
- $\mathbb{R}^X$ : **real space Function**
  - $\mathbb{R}^X = \{h \in \mathbb{R}^X \mid h: X \rightarrow \mathbb{R}\} \Rightarrow \forall x \in X \ h(x) \in \mathbb{R}$
- $S_h = \{(x, h(x))\} \subseteq X \times \hat{Y}$ : **is data set with label**  $h(x)$
- $\hat{Y} \subseteq \mathbb{R}$  then  $\hat{Y}^X \subseteq \mathbb{R}^X$ 
  - $\hat{Y}^X = H = \{h \in \hat{Y}^X \mid h: X \rightarrow \hat{Y}\} \subseteq \hat{Y}^X \Rightarrow \forall x \in X \ h(x) \in \hat{Y}$
  - $l: Y \times \hat{Y} \rightarrow \mathbb{R}^+$  such that  $l(y, \hat{y} \in \hat{Y})$  is loss function between  $y$  and  $\hat{y}$
  - $l_h: X \times Y \rightarrow \mathbb{R}^+$  such that  $l(y, h(x) \in \hat{Y})$  is loss function between  $y$  and  $h(x)$



## 4. Uniform covering numbers for a real-valued function class $H$

Let's assume that  $H$  takes values in some set  $\hat{Y} \subseteq \mathbb{R}$ , so that  $H \subseteq \hat{Y}^X$ .

We will require the **loss function  $l$**  to be bounded. we will assume  $\exists B > 0$  such that:

$$(\forall y \in Y)(\forall \hat{y} \in \hat{Y}) \quad 0 \leq l(y, \hat{y}) \leq B \quad \text{and} \quad l: Y \times \hat{Y} \longrightarrow [0, B]$$

### Definition: The loss function class

We will find it useful to define for any function class  $H \subseteq \hat{Y}^X$  and loss  $l: Y \times \hat{Y} \longrightarrow [0, B]$  the loss function class  $l_H \subseteq [0, B]^{X \times Y}$  given by:

$$l_H = \{l_h: X \times Y \longrightarrow [0, B] \mid l_h(x, y) = l(y, h(x)) \text{ for some } h \in H\}$$

# Uniform convergence in a Real-valued Function class $H$

## Theorem: generalization bound

Let the sets  $Y, \hat{Y} \subseteq \mathbb{R}$ . Let  $H \subseteq \hat{Y}^X$ , and let  $l: Y \times \hat{Y} \rightarrow [0, B]$ .

Let  $D$  be any distribution on  $X \times Y$ .

For any  $\varepsilon > 0$ :

$$\mathbb{P}_{S \sim D^m} \left( \sup_{h \in H} |L_D(h) - L_S(h)| \geq \varepsilon \right) \leq \delta = 4 \mathcal{N}_1 \left( \frac{\varepsilon}{8}, l_H, 2m \right) e^{-m\varepsilon^2/32B^2}$$

# Uniform convergence in a Real-valued Function class $H$

## Lemma: L-Lipschitz loss

Let  $Y, \hat{Y} \subseteq \mathbb{R}$ .

Let  $H \subseteq \hat{Y}^X$ , and let  $l: Y \times \hat{Y} \rightarrow [0, B]$ .

$l$  is Lipschitz in its second argument with Lipschitz constant  $L > 0$ , if and only if:

$$|l(y, \hat{y}_1) - l(y, \hat{y}_2)| \leq L|\hat{y}_1 - \hat{y}_2| \quad \forall y \in Y, \hat{y}_1, \hat{y}_2 \in \hat{Y}$$

Then for any  $m \in \mathbb{N}$

$$\mathcal{N}_1(\varepsilon, l_F, m) \leq \mathcal{N}_1\left(\frac{\varepsilon}{L}, H, m\right)$$

# Uniform convergence in a Real-valued Function class $H$

## Corollary: generalization bound

Let  $Y, \hat{Y} \subseteq \mathbb{R}$ .

Let  $H \subseteq \hat{Y}^X$ , and let  $l: Y \times \hat{Y} \rightarrow [0, B]$  such that  $l$  is Lipschitz in its second argument with Lipschitz constant  $L > 0$ .

Let  $D$  be any distribution on  $X \times Y$ .

For any  $\varepsilon > 0$ :

$$\mathbb{P}_{S \sim D^m} \left( \sup_{h \in H} |L_D(h) - L_S(h)| \geq \varepsilon \right) \leq \delta = 4 \mathcal{N}_1 \left( \frac{\varepsilon}{8L}, H, 2m \right) e^{-m\varepsilon^2/32B^2}$$