

30 August 2021

## CS492(F) Comp. Learning Theory - The PAC Learning Framework (Ch2)

### 1. Overview / Setting

(1) In this chapter, we consider the problem of learning a classifier from examples, and study theoretical tools for analysing such a problem.

(2) A classification problem is described as follows:

①  $\mathcal{X}$  ... space for the input. Usually,  $\mathcal{X} = \mathbb{R}^I$  for some  $I \in \mathbb{N}$ .

②  $\mathcal{Y}$  ... space for the outputs. We assume  $\mathcal{Y} = \{0, 1\}$ .

③ unknown concept  $C$ :  $\mathcal{X} \rightarrow \mathcal{Y}$  from a concept class  $\mathcal{C}$ .

④ unknown distribution  $D$  on  $\mathcal{X}$ .

⑤ Sample  $S = (x_1, x_2, \dots, x_m)$  and  $(y_1, \dots, y_m)$  s.t.

(or training set)  $x_1, \dots, x_m$  are independently drawn from  $D$  and.

$y_i = C(x_i)$  for all  $i \in [m] = \{1, 2, 3, \dots, m\}$ .

⑥ Hypothesis set  $\mathcal{H} \subseteq [\mathcal{X} \rightarrow \mathcal{Y}]$

a good

⑦ Given ①, ②, ⑤, ⑥, find  $h \in \mathcal{H}$  that behaves like  $C$  for the inputs drawn from  $D$ .

(?) To analyse such a problem and its candidate solution rigorously, we need mathematical formalisation of "good" or "behave like", which expresses a desired property of  $h \in \mathcal{H}$ . We want this formalisation also talks about computational resources need to find such a good  $h$ . By computational resources, we mean time and space (the usual suspect in CS), as well as the number of examples in  $S$  needed. The latter is called sample complexity.

and is new in these learning problems.

- (4) The PAC learning framework is such a formalism. It also naturally leads to the idea of generalization bound, which appears repeatedly in the course.
- (5) The purpose of this chapter is to learn key concepts and results in the PAC framework, and to prove generalization bounds for some particular classification problems.

## 2. The PAC learning model:

[ sample (random) ]  
unknown dist.

- (1) Reminder:  $X, Y, \mathcal{C}, \mathcal{H}, \mathcal{D}, S, c \in \mathcal{C}$
- $\xrightarrow{\text{input space}}$   $\xrightarrow{\text{f.o.i.s.}}$   $\xrightarrow{\text{Concept class.}}$   $\xrightarrow{\text{hypo. set.}}$   $\xrightarrow{\text{target}}$   $\xrightarrow{\text{unknown concept.}}$

- (2) The accuracy of a hypo.  $h$  is formalized by the next two concepts, generalisation error (or risk) and empirical error (or empirical risk).

Def 2.1 Given:  $h \in \mathcal{H}, c \in \mathcal{C}, \mathcal{D}$ .

on iid draws from  $\mathcal{D}$ .

The generalisation error or risk of  $h$  is defined by

$$R(h) = \underset{x \sim \mathcal{D}}{\mathbb{P}}[h(x) \neq c(x)] = \underset{x \sim \mathcal{D}}{\mathbb{E}}[\mathbb{1}_{\{h(x) \neq c(x)\}}]$$

error in an ideal setup  
depends on  $c$  and  $\mathcal{D}$ .

indicator for returning 1.  
if the condition holds, 0  
otherwise.

Def 2.2 Given:  $h \in \mathcal{H}, c \in \mathcal{C}, \mathcal{D}, S = (x_1, \dots, x_m) \sim \mathcal{D}^m$

The empirical error or empirical risk of  $h$  is defined by

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq c(x_i)\}} = \underset{x \sim \mathcal{D}}{\mathbb{E}}[\mathbb{1}_{\{h(x) \neq c(x)\}}]$$

error that can be measured in practice.  
depends on  $c, \mathcal{D}$

uniform distribution over:  
 $(x_1, \dots, x_m)$ .

Lemma:  $\mathbb{E}_{S \sim D^m} [\widehat{R}_S(h)] = R(h)$ . That is,  $\widehat{R}_S(h)$  is an unbiased estimate of  $R(h)$ .

$$\begin{aligned} \text{Pr..f. } \mathbb{E}_{S \sim D^m} [\widehat{R}_S(h)] &= \mathbb{E}_{S \sim D^m} \left[ \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq c(x_i)\}} \right] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{x_i \sim D} [\mathbb{1}_{\{h(x_i) \neq c(x_i)\}}] = \frac{1}{m} \sum_{i=1}^m R(h) = R(h) \quad \square. \end{aligned}$$

Let  $m$  be the size of the representation of each  $x \in \mathcal{X}$ .

Def 2.3. [PAC learning]  $\mathcal{X}, \mathcal{Y}$  given.

A concept class  $\mathcal{C}$  is PAC-learnable if  $\exists$  an algo.  $\mathcal{A}$  and a polynomial  $\text{poly}(\cdot, \dots)$  s.t.

$\forall c \in \mathcal{C} \quad \forall D \in \text{Pr}(\mathcal{X}) \quad \forall \varepsilon > 0 \quad \forall \delta > 0 \quad \forall m \in \mathbb{N}$ ,  
 $(D \text{ is a prob. distribution on } \mathcal{X}) \rightarrow$  the size of the representation  
 of  $c$ .

if  $m \geq \text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, n, \text{size}(c))$ , then

$\mathbb{P}_{S \sim D^m} [R(h_S) \leq \varepsilon] \geq 1 - \delta$ .

the result of running the alg.  $\mathcal{A}$  on  
 $\{(x_i, c(x_i)) \mid x_i \in S\}$ .

$\mathcal{C}$  is efficiently PAC-learnable if  $\mathcal{A}$  runs in  $\text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, n, \text{size}(c))$   
 for some polynomial  $\text{poly}$ .

In that case,  $\mathcal{A}$  is called a PAC-learning algo.

\*Note 1: Polynomial sample complexity is the key requirement of PAC learnability.

\*Note 2: Efficient PAC learnability additionally requires polynomial time (or space) complexity.

\*Note 3: No assumption on the distribution  $D$ .

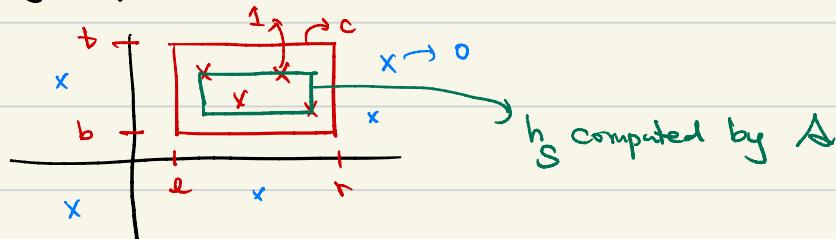
\* Note 4: The definition is applicable usually when  $\exists h$  s.t.  
 $R(h) = 0$ .

\* Note 5:  $\text{size}(c)$  will not play any significant role. We will ignore it  
 (in the course)

(3) Example : Learning axis-aligned rectangles.

$$\mathcal{X} = \mathbb{R}^2, \quad \mathcal{Y} = \{0, 1\}.$$

$$\mathcal{H} = \mathcal{C} = \{(a, b) \mapsto \mathbf{1}_{\{(a, b) \in [l, r] \times [b, t]\}} \mid l \leq r \wedge b \leq t\}.$$



[Q] Prove that  $\mathcal{C}$  is PAC-learnable.

[A]  $A((x_1, y_1), \dots, (x_m, y_m)) =$  the smallest rectangle containing

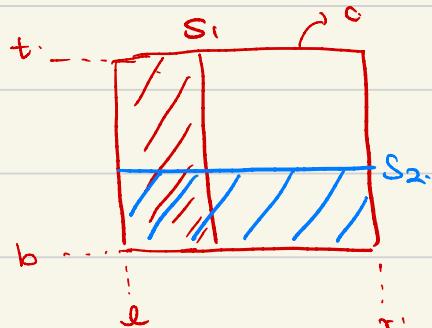
all  $x_i$ 's with  $y_i = 1$ .

$$\text{poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}\right) = \frac{4}{\varepsilon} \log \frac{4}{\delta}.$$

$$\text{Fix } D \in \mathcal{P}(\mathcal{X}), \quad c \in \mathcal{C}, \quad \varepsilon, \delta > 0, \quad m \geq \frac{4}{\varepsilon} \log \frac{4}{\delta}.$$

We will show that

$$\mathbb{P}[R(h_s) > \varepsilon] \leq \delta.$$



Pick  $s_1, s_2, s_3, s_4$  s.t.

$$\mathbb{P}[x \in [l, s_1] \times [b, t]] \geq \varepsilon/4$$

$$\mathbb{P}[x \in [l, s_1] \times [b, t]] \leq \varepsilon/4.$$

$s_2, s_3, s_4$  satisfy similar properties for the other three sides of the rectangle.

$$\text{Let } r_1 = [l, s_1] \times [b, t]$$

$r_2 = [l, r] \times [b, s_2]$  and  $r_3, r_4$  be similar rectangles.

Then,  $R(h_s) > \varepsilon$  implies that  $\bigvee_{i=1}^4 (h_s^{-1}(1) \cap r_i = \emptyset)$ .

$$\begin{aligned}
 \text{Thus, } \Pr_{\text{S} \sim D^m} [R(h_s) > \varepsilon] &\leq \Pr_{\text{S} \sim D^m} \left[ \bigvee_{i=1}^m h_s^{(1)} \cap r_i = \emptyset \right] \\
 &\leq \sum_{i=1}^m \Pr_{\text{S} \sim D^m} [h_s^{(1)} \cap r_i = \emptyset] \leq \sum_{i=1}^m \left(1 - \frac{\varepsilon}{4}\right)^m \leq 4 \exp\left(-\frac{m\varepsilon}{4}\right) \\
 (\because \text{union-bound}) &\quad (\because 1-x \leq e^{-x}) \\
 &\leq 4 \exp\left(-\frac{\varepsilon}{8} \cdot \frac{4}{\varepsilon} \log \frac{4}{\delta}\right) = \delta.
 \end{aligned}$$

□

### 3. PAC learning for finite consistent hypothesis sets.

(1) Assumption :  $\mathcal{H}$  is finite and  $\mathcal{C} \subseteq \mathcal{H}$ .

Thm 2.5: A ... algo. s.t. for any target concept  $c \in \mathcal{H}$  and any  $S = (x_1, \dots, x_m) \in \mathcal{X}^m$ , it returns a consistent  $h_S$ :

$$\widehat{R}_S(h_S) = 0.$$

Then,  $\forall D \in \Pr(\mathcal{X}) \quad \forall \varepsilon, \delta > 0 \quad \forall m > 0$ ,

if  $m \geq \frac{1}{\delta} (\log |\mathcal{H}| + \log \frac{1}{\delta})$ , then  $\Pr_{\text{S} \sim D^m} [R(h_s) \leq \varepsilon] \geq 1 - \delta$ .

□

\*Note 1: Equivalent statement that expresses generalisation bound:

$\forall D \in \Pr(\mathcal{X}) \quad \forall \varepsilon, \delta > 0 \quad \exists m > 0$

$\Pr_{\text{S} \sim D^m} [R(h_s) \leq \frac{1}{m} (\log |\mathcal{H}| + \log \frac{1}{\delta})] \geq 1 - \delta$ .

if  $\log |\mathcal{H}|$  is polynomial in  $m$  or  
(i.e., the size of input)

\*Note 2: Thm 2.5 implies PAC-learnability, although it doesn't imply efficient PAC-learnability.

algo.-independent  
uniform bound.

Proof: We will show that

$\Pr_{\text{S} \sim D^m} [\exists h \in \mathcal{H} \text{ st. } \widehat{R}_S(h) = 0 \wedge R(h) > \varepsilon] \leq \delta.$

Since  $A$  returns a consistent hypo.  $h_S$ , the above inequality implies the thm..

$$\begin{aligned}
& \mathbb{P} [\exists h \in \mathcal{H} \text{ s.t. } \widehat{R}_S(h) = 0 \wedge R(h) > \varepsilon] \\
& \leq \sum_{\substack{h \in \mathcal{H} \\ R(h) > \varepsilon}} \mathbb{P} [\widehat{R}_S(h) = 0] \quad \text{by union bound} \\
& = \sum_{\substack{h \in \mathcal{H} \\ R(h) > \varepsilon}} \prod_{i=1}^m \mathbb{P}[h(x_i) = c(x_i)] \leq \sum_{h \in \mathcal{H}} (1-\varepsilon)^m = |\mathcal{H}| (1-\varepsilon)^m \leq |\mathcal{H}| e^{-\varepsilon m} \\
& = |\mathcal{H}| e^{-\varepsilon \times \frac{1}{\varepsilon} \log |\mathcal{H}| + \log \frac{1}{\varepsilon}} = g. \\
& \xrightarrow{\text{assumption on } m.} \quad \text{D.}
\end{aligned}$$

uniform bound. exponential decay  
due to repeated trials.

\* Note 1 :  $\frac{1}{m} \log |\mathcal{H}|$  comes from the use of uniform bound,  
 and  $\frac{1}{m} \log \frac{1}{\varepsilon}$  comes from the estimation error for  
 a single  $h$ .

\* Note 2 :  $O(\frac{1}{m})$  convergence rate. In the generalisation bound.  
 This is a good news. In a more difficult inconsistent case,  
 we will get  $O(\frac{1}{\sqrt{m}})$  convergence rate.

\* Note 3 :  $\log |\mathcal{H}|$  corresponds to the # of bits needed to  
 encode each hypothesis in  $\mathcal{H}$ .

## (2) Examples / Applications.

### ① Conjunctions of boolean literals.

- $m$  boolean variables  $v_1, \dots, v_n$
- each example specifies the values of those variables.  $\mathcal{X} = \mathbb{Z}_2^n$   
 $x = (v_1=1, v_2=1, v_3=0, \dots, v_n=1)$  (denoted by  $(110 \dots 1)$ )
- $\mathcal{C} = \mathcal{L} = \{v_1 \wedge \bar{v}_2, v_1 \wedge v_2 \wedge \bar{v}_3, \dots\} = \text{the set of conjunctions of literals}$
- $\log |\mathcal{L}| = \log 3^m = m \log 3$ .

## ② Universal Concept class

- $$- \quad \mathcal{U}_n = [\mathbb{X} \rightarrow \mathbb{Y}] \ , \quad \quad \mathbb{X} = 2^n \ , \ \mathbb{Y} = \{0,1\}.$$

i.e., the set of all subsets of  $\mathbb{N}$ ).

$$ff \geq 2n.$$

- Sample complexity bound from Thm 2.5

$$\begin{aligned}
 m &\geq \frac{1}{\varepsilon} (\log |f| + \log \frac{1}{\delta}) \geq \frac{1}{\varepsilon} (\log |h_n| + \log \frac{1}{\delta}) \\
 &= \frac{1}{\varepsilon} (\log 2^{2^n} + \log \frac{1}{\delta}) \\
 &= \frac{1}{\varepsilon} (\underbrace{2^n \log 2}_{\text{exponential in } n} + \log \frac{1}{\delta})
 \end{aligned}$$

- In fact, not PAC-learnable.

### ③ k-term DNF formulas.

- $\text{Sf} = \mathcal{C} = \{ \bigvee_{j \in [u]} C_j \mid C_j \text{ is a conjunction of literals}$   
 $\sum_{i=1,2,\dots,u} \{ \dots \}$  and  $u \leq k \}$ .

$$|\mathcal{F}| = |\mathcal{C}| = 3^{n \times 0} + 3^{n \times 1} + 3^{n \times 2} + \dots + 3^{n \times k} \leq (k+1) 3^{n \times k}$$

- Sample complexity bound from Thm 2.5

$$\frac{1}{\epsilon} (\log |f| + \log \frac{1}{\delta}) \leq \frac{1}{\epsilon} (m \times k \times \log 3 + \log(k+1) + \log \frac{1}{\delta}). \dots \text{poly.}$$

PAC-learnable.

- But not efficiently PAC-learnable if RP  $\neq$  NP.

#### ④ k-CNF

$\mathcal{C} = \text{ff} = \{ \bigwedge_{j=1}^u C_j \mid C_i \neq C_j \text{ for } i \neq j, C_j \text{ is a disjunction of at most } k \text{ literals} \}$

variable  $v$  or its negation

- Reduction to the case ① by replacing all disjunctions of at most  $k$  literals,  $(l_1 \vee l_2 \vee \dots \vee l_k)$ , with new variables  $v_{l_1}, v_{l_2}, \dots, v_{l_k}$ . Then, the # of these new variables  $\leq (2n+1)^k$ .
- The sample complexity bound from Thm 2.4 is  $m \geq \frac{1}{\epsilon} ((2n+1)^k \log 3 + \log \frac{1}{\delta})$ .
- PAC-learnable. In fact, efficiently PAC-learnable.

#### 4. Guarantee for finite inconsistent-hypothesis sets.

- We now consider the case that the hypo. set  $\text{ff}$  is finite but  $\mathcal{C} \neq \text{ff}$ , more accurately,  $\exists c \in \mathcal{C}$  and  $D \in \text{Pf}(x)$  s.t.  $R(h) \rightarrow 0$  for any  $h \in \text{ff}$  when  $R$  is defined with  $c$  and  $D$ .
- In such a case, PAC-learning is not possible usually, because sometime we cannot reduce  $R(h)$  enough, no matter how many examples we use in  $S = (x_1, \dots, x_m)$ .
- We instead try to find a bound on generalization error which has the form: high-prob. uniform.

$$\Pr[\forall h \in \text{ff}, R(h) \leq \hat{R}_{S,h} + f(m, \delta, |\text{ff}|)] \geq 1 - \delta.$$

Such a bound lets us talk about the quality of a learnt classifier measured by  $R(h_s)$  using the estimated  $\hat{R}_{S(h_s)}$ .

Thm 2.13. If  $f$  is finite. uniform bound.

$$\Rightarrow \forall c \in \mathcal{C} \quad \forall D \in \mathcal{D} \quad \Pr[X] \quad \forall \epsilon > 0 \quad \forall \delta > 0.$$

$$\Pr_{S \sim D^m} [\forall h \in f. \quad R(h) \leq \hat{R}_{S(h)} + \sqrt{\frac{\log |f| + \log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

$O(\sqrt{\frac{1}{m}})$ , instead of  $O(\frac{1}{\sqrt{m}})$

Equivalent statement based on sample complexity bound.

$$\forall c \quad \forall D \quad \forall \epsilon, \delta > 0 \quad \forall m \geq \frac{1}{\epsilon^2} (\log |f| + \log \frac{1}{\delta})$$

$$\Pr [\forall h \in f. \quad R(h) \leq \hat{R}_{S(h)} + \epsilon] \geq 1 - \delta.$$

(\*) Proof of thm 2.13:

① Hoeffding's inequality. ... summing ind. random vars leads to concentration around the mean.

Thm D.2  $x_1, \dots, x_m$  ... independent random variables. s.t.

$x_i \in [a_i, b_i]$ . for all  $i \in [m]$ .

$$S_m = \sum_{i=1}^m x_i$$

$$\epsilon > 0$$

$\Rightarrow$

$$\Pr [S_m - \mathbb{E}[S_m] \geq \epsilon] \leq \exp \left( -2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2 \right)$$

$$\text{and } \Pr [S_m - \mathbb{E}[S_m] \leq -\epsilon] \leq \exp \left( -2\epsilon^2 / \sum_{i=1}^m (b_i - a_i)^2 \right)$$

Cor 2.11 Let  $h \in f$ , and  $\delta > 0$ . Fix  $c, D, m$ . Then,

$$\Pr_{S \sim D^m} [R(h) \leq \hat{R}_{S(h)} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

Prf. It suffices to show that

$$\mathbb{P} [ R(h) \geq \hat{R}_S(h) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} ] \leq \delta.$$

$$\begin{aligned} &= \mathbb{P} \left[ \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq c(x_i)\}} - \mathbb{E} \left[ \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq c(x_i)\}} \right] \leq \sqrt{\frac{m \log \frac{1}{\delta}}{2}} \right] \\ &\stackrel{\text{Hoeffding}}{\leq} \exp \left( \frac{-2 \times \left( \sqrt{\frac{m \log \frac{1}{\delta}}{2}} \right)^2}{\sum_{i=1}^m (1-\alpha)^2} \right) = \exp \left( \frac{-2 \times \frac{m}{2} \log \frac{1}{\delta}}{m} \right) \\ &= \delta. \quad \square. \end{aligned}$$

Equivalent version of ..

for all  $\epsilon > 0$ ,

$$\mathbb{P} [ R(h) \geq \hat{R}_S(h) + \epsilon ] \leq \exp(-2m\epsilon^2)$$

② Proof of Thm 2.13.

$$\text{Let } \epsilon = \sqrt{\frac{\log |\mathcal{H}| + \log \frac{1}{\delta}}{2m}}.$$

We will show that

$$\begin{aligned} &\mathbb{P} [ \exists h \in \mathcal{H} \text{ s.t. } R(h) \geq \hat{R}_S(h) + \epsilon ] \leq \delta. \\ &\leq \sum_{h \in \mathcal{H}} \mathbb{P} [ R(h) \geq \hat{R}_S(h) + \epsilon ]. \\ &\leq \sum_{h \in \mathcal{H}} \exp(-2m\epsilon^2) = |\mathcal{H}| \exp \left( -2m \times \left( \frac{\log |\mathcal{H}| + \frac{1}{\delta}}{2m} \right) \right) \\ &= \delta. \quad \square. \end{aligned}$$

(c5) Notes on Thm 2.13.

$$① R(h) \leq \hat{R}_S(h) + O \left( \sqrt{\frac{\log |\mathcal{H}|}{m}} \right)$$

Contrast it with the generalization bound of Thm 2.5

$$R(h) \leq O \left( \frac{\log_2 |\mathcal{H}|}{m} \right).$$

Thus, worse bound.

② Simpler  $f$  and larger  $m$  tighten the bound.

related to Occam's razor principle. ... the simplest explanation is best.

## 5. Stochastic scenarios

(1) A natural generalisation of the setting of a distribution  $D$  on the inputs  $\mathcal{X}$  and an unknown concept  $c$  is to assume simply a distribution  $D'$  on the input-output pairs. Note that the previous setting is a special case of this stochastic general setting where  $D'(x,y)$  is defined to be  $D(x) \times \mathbb{1}_{\{y=c(x)\}}$ . (More formally,  $D'$  is the push-forward measure of  $D$  along the function  $x \mapsto (x, c(x))$ .)

(2) In this new stochastic setting, we cannot avoid a certain level of error, and this is not compatible with the PAC framework that usually works when we can minimise error as much as we can. So, we have to modify PAC.

(3) New setup and agnostic PAC learning.

$$D \in \Pr(\mathcal{X} \times \mathcal{Y})$$

$$S = ((x_1, y_1), \dots, (x_m, y_m)) \sim D^m$$

$$R(h) = \mathbb{P}_{(x,y) \sim D} [h(x) \neq y] = \mathbb{E}_{(x,y) \sim D} [\mathbb{1}_{\{h(x) \neq y\}}].$$

$$\widehat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq y_i\}}.$$

 Def 2.14 If  $\mathcal{H}$  ... hpo. set.  $\Delta$  ... a learning algorithm.

$\Delta$  is an agnostic PAC-learning algo. if  $\exists$  a polynomial  $\text{poly}(\cdot, \cdot)$  s.t.  $\forall D \in \mathbb{P}(\mathcal{X} \times \mathcal{Y}) \quad \forall \varepsilon, \delta > 0 \quad \exists m \in \mathbb{N}$

if  $m \geq \text{poly}\left(\frac{1}{\varepsilon}, \frac{1}{\delta}, \overset{\text{size of each input}}{n}\right)$ , then

$$\Pr_{S \sim D^m} [\text{RCh}_S - \min_{h \in \mathcal{H}} \text{R}(h) \leq \varepsilon] \geq 1 - \delta.$$

$\mathcal{X}$  is finite and error that cannot be avoided.

(\*) Suppose that  $\mathcal{H} = [\mathcal{X} \rightarrow \mathcal{Y}]$ , the set of all functions.

What is the solution of the following opt. pb?

$$\min_{h \in \mathcal{H}} \text{R}(h) = \min_{h \in \mathcal{H}} \mathbb{E}[ \mathbb{1}_{f(h(x)) \neq y} ].$$

- Optimal solution  $h_{\text{Bayes}}$  (called Bayes hypothesis, Bayes classifier)

$$h_{\text{Bayes}}(x) = \underset{y \in \{0, 1\}}{\operatorname{argmax}} \Pr[y | x].$$

- Optimal value  $R^* = R(h_{\text{Bayes}}) = \mathbb{E}[\text{noise}(x)]$  Bayes error.

where  $\text{noise}(x) = \min \{ \Pr[1|x], \Pr[0|x] \}$ .

called noise at  $x$ .

- For each  $D \in \mathbb{P}(\mathcal{X} \times \mathcal{Y})$ , the Bayes error describes the difficulty of the learning problem wrt.  $D$ .