

# FOUNDATIONS OF STATISTICAL MACHINE LEARNING

## CLASSICAL LEARNING MODEL

Domain Set  $\rightarrow$  Set of **objects** for labelling ( $X$ )  
 $\Downarrow$  Depicted as feature vectors  
(input to the algorithm)

Label Set  $\rightarrow$  Set of possible labels ( $Y$ )

Training Data  $\rightarrow S \subseteq X \times Y$   
Feature vectors + Labels.

Learner's Output  $\rightarrow$  Hypothesis function  $h: X \rightarrow Y$   
Prediction Algorithm  $\rightarrow A$

Distribution  $\rightarrow$  Distribution  $\mathcal{D}$  exists on domain  $X$   
For labels  $\rightarrow$  There exists correct labelling  
function ( $f$ )

Loss Function  $\rightarrow$

$$L_{\mathcal{D}, f}(h) = \Pr_{x \in \mathcal{D}} [h(x) \neq f(x)]$$

# EMPIRICAL RISK MINIMISATION

$f \rightarrow$  unknown  $D \rightarrow$  unknown

★ Minimising error over training data

Learner has  $S$  according to  $D$   
SAMPLE DISTRIBUTION

$$S = (x_1, y_1), (x_2, y_2) \dots (x_m, y_m)$$

where  $x_i \sim D$

$$y_i = f(x_i)$$

$$L_S = \frac{|\{i \in [m] \mid h(x_i) \neq f(x_i)\}|}{m}$$

## OVERFITTING :-

The distribution  $D$  may confuse a learner to learn a **wrong rule**.  $\Rightarrow$  Because  $S$  is a part of  $D$  and does not represent  $D$  exhaustively

This rule will work fine on training data

ERM  $\rightarrow$  Find  $h$  that minimizes  $L_S$

★ Finding  $h$  may not be an efficient way because all possible functions may be very large.

PROBLEM: Aiming for  $L_S(h) = 0$  may lead to overfitting.

★ USE ADDITIONAL INFO:  $ERM_{\mathcal{H}}$   $\mathcal{H} \rightarrow$  Hypothesis Class

$$ERM_{\mathcal{H}}(S) = \argmin_{h \in \mathcal{H}} L_S(h)$$

PRIOR KNOWLEDGE

★ A restrictive  $\mathcal{H}$  may make learning the problem easy.

★ If  $\mathcal{H}$  is restricted, it can lead to large INDUCTIVE BIAS

# Over which hypothesis classes  $\mathcal{H}$  will  $ERM_{\mathcal{H}}$  not lead to overfitting??

$$h_S = \argmin_{h \in \mathcal{H}} L_S(h)$$

Result of applying  $ERM_{\mathcal{H}}$

non-intuitive

★ REALIZABILITY ASSUMPTION:  $f \in \mathcal{H}$

## ACCURACY PARAMETER:

$L_{D,f}(h_S) \Rightarrow$  Probability of choosing an example  $x \sim D$  such that  $h_S(x) \neq f(x)$

⊛ If  $L_{D,f}(h_S) < \epsilon$ , then learner is approximately correct

$1 - \epsilon \Rightarrow$  Accuracy Parameter

THEOREM: Let  $\mathcal{H}$  be a finite hypothesis class of functions from domain  $X$  to range  $\{0,1\}$ . Then for every  $f \in \mathcal{H}$  and every  $D$  over  $X$ ,

$$\Pr_{S \sim D^m} [L_{D,f}(\text{ERM}_{\mathcal{H}}(S)) > \epsilon] \leq |\mathcal{H}| \cdot (1-\epsilon)^m$$

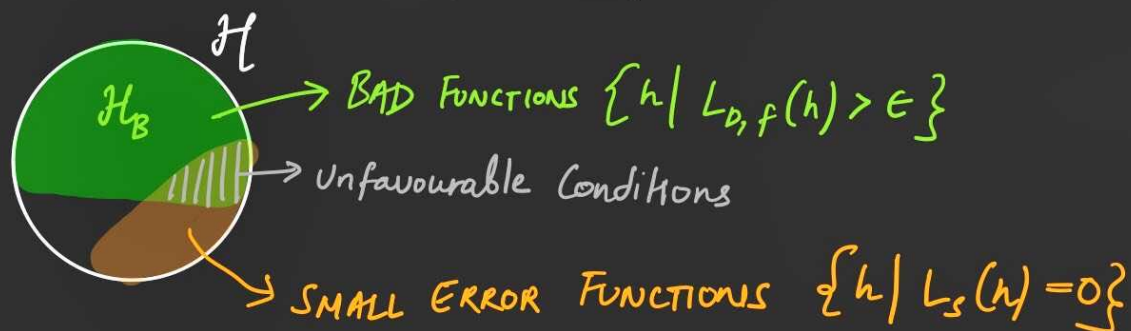
$\downarrow$  Sample size  $m$        $\downarrow$  Size of  $\mathcal{H}$        $\downarrow$  Probability goes down exponentially

★ Theorem holds for every  $m$  and every  $\epsilon$

★ As long as  $e^{-\epsilon m}$  dominates  $|\mathcal{H}|$  we can contain error in  $L_{D,f}(h_S)$

★ Therefore, if  $|\mathcal{H}|$  is bounded or finite, then we can bound  $|\mathcal{H}| \cdot (1-\epsilon)^m$  by changing sample size  $m$ .

PROOF:



GOAL: To bound  $\Pr[\exists h \in \mathcal{H}_B \text{ such that } L_S(h) = 0]$

$$\Pr[\exists h \in \mathcal{H}_B \text{ and } L_S(h) = 0] \leq |\mathcal{H}_B| \Pr_{S \sim D^m} [L_S(h) = 0] \xrightarrow{h \in \mathcal{H}_B}$$

Assumption about  $\mathcal{H}_B \Rightarrow L_{D,f}(h) > \epsilon$

Thus  $\Pr_{S \sim D^m} [L_S(h) = 0 \text{ for } h \in \mathcal{H}_B] \leq (1-\epsilon)^m$

Finally  $\Pr[\exists h \in \mathcal{H}_B \text{ and } L_S(h) = 0] \leq |\mathcal{H}_B| (1-\epsilon)^m \leq |\mathcal{H}| (1-\epsilon)^m$



## PROBABILISTICALLY APPROXIMATELY CORRECT

A class  $\mathcal{H}$  is said to be PAC learnable, if there is a function  $m_{\mathcal{H}}: (0,1) \times (0,1) \rightarrow \mathbb{N}$  and learning algorithm  $A$  such that for every  $\epsilon, \delta \in (0,1)$ , for every Distribution  $\mathcal{D}$  and for every  $f \in \mathcal{H}$

$$\Pr_{S \sim \mathcal{D}^m} [L_{\mathcal{D},f}(A(S)) \leq \epsilon] > 1 - \delta$$

for any  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$

$\epsilon \rightarrow$  Accuracy Parameter  
 $1 - \delta \rightarrow$  Confidence Parameter

What Algorithm to use?  $\Rightarrow \text{ERM}_h$

Since  $|\mathcal{H}|(1-\epsilon)^m \leq |\mathcal{H}|e^{-\epsilon m} \leq \delta$

$$m_{\mathcal{H}}(\epsilon, \delta) = \frac{\ln(|\mathcal{H}|) + \ln(1/\delta)}{\epsilon}$$

STRENGTH: Valid for any distribution

WEAKNESS: Assumption that  $f \in \mathcal{H}$   $\rightarrow$

# How to remove this assumption??

used when we assumed  $L_S(h) = 0$   
if  $h = \text{ERM}_{\mathcal{H}}(S)$

\* Agnostic PAC learning

## Removing realisability assumption

$D \rightarrow$  Distribution over  $X \times Y$

$D_x \rightarrow$  Distribution over unlabelled domain points  $\rightarrow$  Marginal

$D((x, y) | x) \rightarrow$  conditional probability of labels for each domain point.

TRUE ERROR / RISK : The likelihood of  $h$  making an error when labelled points are drawn as per  $D$

$$L_D(h) = \Pr_{(x, y) \sim D} [h(x) \neq y]$$

EMPIRICAL RISK :

$$L_S(h) = \frac{|\{i \in [m] \mid h(x_i) \neq y_i\}|}{m}$$

## BAYES OPTIMAL PREDICTOR

$$f_D(x) = \begin{cases} 1 & \text{if } \Pr[y=1|x] \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

Prove that: For any classifier  $g$ ,  $L_D(f_D) \leq L_D(g)$

PROBLEM : Learner does not know  $D$

# AGNOSTIC PAC LEARNING

A class  $\mathcal{H}$  is said to be agnostic PAC learnable, if there is a function  $m_{\mathcal{H}} : (0,1) \times (0,1) \rightarrow \mathbb{N}$  and a learning algorithm  $A$  such that for every distribution  $D$  on  $X \times Y$ ,

$$\Pr_{S \sim D^m} [L_D(A(S)) \leq \min_{h \in \mathcal{H}} \{L_D(h)\} + \epsilon] > 1 - \delta$$

for any  $m \geq m_{\mathcal{H}}(\epsilon, \delta)$

→ cannot guarantee 0 error  
This is best possible function

STRENGTH: Realizability assumption not necessary

\* If  $H_1 \subseteq H_2$ ,  $H_1$  is easier to learn.

WEAKNESS: The error is relative  
\* with respect to  $\min_{h \in \mathcal{H}} \{L_D(h)\}$

## OTHER LEARNING TASKS:

- ① Multiclass Prediction
- ② Real valued Prediction (Regression)

→ Different notion of loss

### General Setup for Learning:

① Binary Label Prediction

$$\text{loss} = \begin{cases} 1 & h(x) \neq y \\ 0 & h(x) = y \end{cases}$$

② Multiclass Prediction

$$\text{loss} = \begin{cases} 1 & h(x) \neq y \\ 0 & h(x) = y \end{cases}$$

③ Regression

$$l(h, (x, t)) = (h(x) - t)^2$$

④ Classification

$$l((c_1 \dots c_k), z) = \min_{1 \leq i \leq k} \|c_i - z\|^2$$

# LEARNING PARADIGM: UNIFORM CONVERGENCE

DEFINITION: A sample  $S = (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  is  $\epsilon$ -representative of a class  $H$  with respect to the distribution  $D$  if  $\forall h \in H \quad |L_S(h) - L_D(h)| \leq \epsilon$  where  $L_S(h) = \frac{1}{n} \sum_{(x, y) \in S} \ell(h, (x, y))$

\* Note that if  $S$  is indeed representative of  $H$  with respect to  $D$  then  $ERM_H$  is a good learning strategy

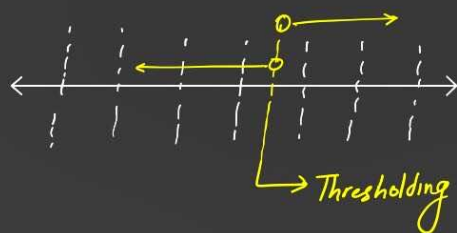
CLAIM: If  $S$  is  $\epsilon$ -representative of  $H$  with respect to  $D$  then for any  $ERM_H$  function  $h_S$

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

PROOF:  $L_D(h_S) \leq L_S(h_S) + \epsilon \leq \min_{h \in H} [L_S(h) + \epsilon] \leq \min_{h \in H} [L_D(h)] + 2\epsilon$

\* Practical approach to infinite classes  $\rightarrow$  Discretization

EXAMPLE  $\rightarrow H_\kappa^{thr} = \{h_r: h_r(x) = \begin{cases} 0 & x \leq r \\ 1 & x > r \end{cases}, r \in \{0, \frac{1}{\kappa}, \frac{2}{\kappa}, \dots, 1\}\}$



\* In theory,  $H_\kappa^{thr}$  may not be a good approximation of  $H^{thr}$ .

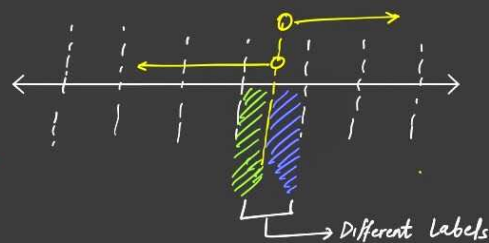
$$\min_{h \in H^{thr}} L_D(h) \ll \min_{h \in H_\kappa^{thr}} L_D(h)$$

EXAMPLE

$$D(\{(\frac{3}{4}, 0)\}) = 0.5$$

$$D(\{(\frac{3}{4}, 1)\}) = 0.5$$

$\rightarrow$  Put everything on 2 point



NO FREE LUNCH THEOREM:

Let  $X$  be a domain of size  $n$ .

Let  $H_n^{all}$  be the set of all possible labellings.

$$H_n^{all} = \{h: X \rightarrow [0, 1]\} \quad |H_n^{all}| = 2^n$$

NFL proves that  $M_{H_n^{all}}(\frac{1}{2}, \frac{1}{7}) \geq \frac{n}{2}$

Therefore, the class of labelling functions over an infinite domain is not PAC-learnable.



## TRADE OFF :

$$\epsilon_{app} = \min_{h \in H} L_D(h) \quad \epsilon_{est} = L_D(h_s) - \epsilon_{app}$$

→ By Free Lunch Theorem this is large.

## THEOREM :

Let  $A$  be a learning algorithm for the task of binary classification over the domain set  $X$ . Let  $m < |X|/2$ . Then  $\exists D$  over  $X \times \{0, 1\}$  such that

$$\exists f: X \rightarrow \{0, 1\} \text{ such that } L_D(f) = 0$$

$$\Pr_{S \sim D^m} \{L_D(A(S)) \geq 1/8\} \geq 1/7$$

## MOTIVATION FOR VC DIMENSION :

→ Comes from proof of no-free-lunch theorem.

### ① NOTION OF SHATTERING A SET

A set of points  $S \subseteq X$  is said to be shattered by  $H$  for every possible labellings of points in  $S$  as 0-1, there is a function  $h \in H$  that realises this labelling.

### EXAMPLE 1 :

$H = \{[0, b] \mid b \in \mathbb{R}^{>0}\}$  \*  $H$  can shatter any set of size 1. **EXCEPTION:  $x=0$  case.**  
\*  $H$  cannot shatter any sets of size 2.

### EXAMPLE 2 :

$H = \{[a, b] \mid a, b \in \mathbb{R} \text{ and } a \leq b\}$  \*  $H$  can shatter any set of size 2  
\*  $H$  cannot shatter any set of size 3

### EXAMPLE 3 :

$H$  is a set of linear functions \*  $H$  can shatter any set of size 2  
\*\*  $H$  can shatter **SOME** sets of size 3

SHATTERED

UNSHATTERED

## VAPNIK - CHERVONENKIS DIMENSION :

Let  $H$  be a hypothesis class over functions from  $X$  to  $\{0, 1\}$ .

The VC-Dimension of  $H$  is the size of the largest finite subset of  $X$  such that it is shattered by  $H$ .

① If  $\exists$  a subset of size  $d$  such that it can be shattered by  $H$ , then VC dimension of  $H$  is at least  $d$ .

② If no subset of size  $d$  is shattered by  $H$ , then VC dimension of  $H$  is less than  $d$ .