

1. Notations: (Supervised Setting). (A Formal Learning Model.) $S = (x_1, y_1), \dots, (x_m, y_m)$ Training set. X : domain set. Y : label set. \mathcal{H} : hypothesis class, a fixed subset of $\{\text{all functions } h: X \rightarrow Y\}$. \mathcal{D} : a dist. on X (to be revealed later).Assumption:① $x_1, \dots, x_m \in S$ is generated iid w/ \mathcal{D} .② y_1, \dots, y_m are labelled by a certain $f: X \rightarrow Y$. labelling function
i.e. $y_i = f(x_i)$, $\forall i$. (to be revealed.) $\ell: \mathcal{H} \times (X \times Y) \rightarrow \mathbb{R}_+$, a loss function, measuring how "preferred" a hypothesis is.For $h \in \mathcal{H}$, $L_S(h) := \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, y_i)) = \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, f(x_i)))$, called training error, empirical error, or empirical risk. $L_{\mathcal{D}}(h) := E_{x \sim \mathcal{D}} [\ell(h, (x, f(x)))]$, called generalization error, risk, or true error.

2.

• An algorithm A is a process of generating an output hypothesis $A(S)$ given the input $S = (x_1, y_1), \dots, (x_m, y_m)$.
 \mathcal{H} • ERM (= Empirical Risk Minimization) is the algorithm s.t. $A(S) \in \argmin_{h \in \mathcal{H}} L_S(h)$.Def Given \mathcal{D} over X . \mathcal{D} is called realizable by \mathcal{H} if $\exists h^* \in \mathcal{H}$ s.t. $L_{\mathcal{D}}(h^*) = 0$.Def (PAC = probably approximately correct).Given \mathcal{H} and $\ell: \mathcal{H} \times (X \times Y) \rightarrow \mathbb{R}_+$.The hypothesis class \mathcal{H} is called PAC learnable if

$\exists m_H: (0,1)^2 \rightarrow \mathbb{N}$ and algorithm A s.t.

\forall dist. \mathcal{D} on \mathcal{X} , $\epsilon, \delta \in (0,1)$ and labelling function $f: \mathcal{X} \rightarrow \mathcal{Y}$,

if the realizable assumption holds, then

running A on $S \sim \mathcal{D}^m$ w/ $m \geq m_H(\epsilon, \delta)$, we have

$$L_{\mathcal{D}}(A(S)) \leq \epsilon \text{ w/ prob. } \geq 1 - \delta. (\text{over } S \sim \mathcal{D}^m).$$

Rank:

① ϵ : accuracy parameter, δ : confidence parameter
(approximately correct) (probably).

② Minimum among such m_H is called the sample complexity.

Issues:

① Realizability assumption is NOT practical (i.e. $f \notin \mathcal{H}$ often happens).

② Labelling function is not realistic (i.e. it's possible to have $x_i = x_j$ but $y_i \neq y_j$).

Sol'n to ②:

Instead of a dist. \mathcal{D} on \mathcal{X} and labelled by f , we can directly consider a dist. \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$. Denote $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$ for simplicity.

To also solve issue ①, we modify PAC learnability as follows:

Def Given \mathcal{H} and $\ell: \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$.

\mathcal{H} is called agnostic PAC learnable if $\exists m_H: (0,1)^2 \rightarrow \mathbb{N}$ and algorithm A

s.t. \forall dist. \mathcal{D} on \mathcal{Z} , $\epsilon, \delta \in (0,1)$,

if we run A on $S \sim \mathcal{D}^m$ w/ $m \geq m_H(\epsilon, \delta)$, we have

$$L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon, \text{ w/ prob. } \geq 1 - \delta (\text{over } S \sim \mathcal{D}^m),$$

where $\mathcal{L}_{\mathcal{D}}(h) := \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$, $\forall h \in \mathcal{H}$.

3. VC dimension and Fundamental Thm of Statistical Learning.

Q: Which classes \mathcal{H} are (agnostic) PAC learnable?

Def

\mathcal{H} : a hypothesis class of fun. from \mathcal{X} to $\{0,1\}$.

$$C = \{c_1, \dots, c_m\} \subseteq \mathcal{X}.$$

① $\mathcal{H}_C := \{f|_C : C \rightarrow \{0,1\} \mid f \in \mathcal{H}\}$, called the restriction of \mathcal{H} to C .

② \mathcal{H} is said to shatter C if $\mathcal{H}_C = \{\text{all fun. from } C \text{ to } \{0,1\}\}$.

i.e. $|\mathcal{H}_C| = 2^{|C|}$.

③ The VC dimension of C , denoted $\text{VCdim}(C)$, is the maximal size of $C \subseteq \mathcal{X}$ that can be shattered by \mathcal{H} .

Rank:

VC = Vapnik-Chervonenkis

Thm (F.T. of S.L. - Qualitative)

\mathcal{H} : a hypothesis class of fun. from \mathcal{X} to $\{0,1\}$, w/ 0-1 loss.

Then T.F.A.E.

(1) \mathcal{H} : agnostic PAC learnable.

(2) \mathcal{H} : PAC learnable.

(3) ERM is a successful agnostic PAC learner for \mathcal{H} .

(4) " " " " " " PAC learner for \mathcal{H} .

(5) $\text{VCdim}(\mathcal{H}) < \infty$.

Thm (F.T. of S.L. - Quantitative)

\mathcal{H} : hypo. class of fun. from \mathcal{X} to $\{0,1\}$ w/ 0-1 loss, $\text{VCdim}(\mathcal{H}) = d < \infty$.

Then \exists abs. const. C_1 and C_2 s.t.

(1) \mathcal{H} : agnostic PAC learnable w/ $C_1 \cdot \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log(1/\delta)}{\epsilon^2}$.

(2) \mathcal{H} : PAC learnable w/ $C_1 \cdot \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \cdot \frac{d + \log(1/\delta)}{\epsilon}$.