

PROPOSITION 3.1 (The Realizability Assumption). Given a set \mathcal{H} of N hypotheses h_1, \dots, h_N and a loss function ℓ , show that the algorithm outputs the hypothesis \hat{h} with probability 1 over random samples S after the number of samples m is sampled according to P and are related by $\ell(\hat{h}, S) \leq \min_{h \in \mathcal{H}} \ell(h, S)$.

More generally, we use the **Probably Approximately Correct (PAC)** learning.

DEFINITION 3.2 (PAC Learning). A hypothesis class \mathcal{H} is PAC learnable if there exists a learning algorithm A such that for any loss function ℓ and any distribution P over \mathcal{X} , for any $\epsilon > 0$ and $\delta > 0$, there exists a sample size m such that for any sample S of size m , the algorithm outputs a hypothesis \hat{h} such that $\ell(\hat{h}, S) \leq \min_{h \in \mathcal{H}} \ell(h, S) + \epsilon$ with probability at least $1 - \delta$.

THEOREM 3.3 (PAC Learning). A hypothesis class \mathcal{H} is PAC learnable if and only if $\text{VCdim}(\mathcal{H}) < \infty$. If $\text{VCdim}(\mathcal{H}) = d$, then the sample size m must be at least $\frac{1}{\epsilon^2} \log \frac{1}{\delta}$ for any loss function ℓ and any distribution P over \mathcal{X} .

DEFINITION 3.3 (Approximate PAC Learnability). A hypothesis class \mathcal{H} is approximately PAC learnable if there exists a learning algorithm A such that for any loss function ℓ and any distribution P over \mathcal{X} , for any $\epsilon > 0$ and $\delta > 0$, there exists a sample size m such that for any sample S of size m , the algorithm outputs a hypothesis \hat{h} such that $\ell(\hat{h}, S) \leq \min_{h \in \mathcal{H}} \ell(h, S) + \epsilon$ with probability at least $1 - \delta$.

THEOREM 6.1 (The Fundamental Theorem of Statistical Learning). Let \mathcal{H} be a hypothesis class of functions from a domain \mathcal{X} to $\{0, 1\}$ and let the loss function be the 0-1 loss. Then, the following are equivalent:
1. \mathcal{H} has the uniform convergence property.
2. Any ERM rule is a successful approximate PAC learner for \mathcal{H} .
3. \mathcal{H} is approximately PAC learnable.
4. \mathcal{H} is PAC learnable.
5. Any ERM rule is a successful PAC learner for \mathcal{H} .
6. \mathcal{H} has a finite VC dimension.

DEFINITION 4.1 (representative sample). A training set S is called a representative sample for a domain \mathcal{X} , hypothesis class \mathcal{H} , loss function ℓ , and distribution P if $\ell(\hat{h}, S) - \ell(\hat{h}, P) \leq \epsilon$.

DEFINITION 4.2 (Uniform Convergence). We say that a hypothesis class \mathcal{H} has the uniform convergence property (w.r.t. a domain \mathcal{X} and a loss function ℓ) if there exists a function $\phi(\epsilon, \delta)$ such that for every $\epsilon > 0$ and $\delta > 0$, and for every probability distribution P over \mathcal{X} , if S is a sample of size $m \geq \phi(\epsilon, \delta)$, then with probability at least $1 - \delta$, S is a representative sample.

We have already seen that $1 \rightarrow 2$ in Chapter 4. The implications $2 \rightarrow 3$ and $3 \rightarrow 4$ are trivial and so is $2 \rightarrow 5$. The implications $4 \rightarrow 6$ and $5 \rightarrow 6$ follow from the No-Free-Lunch theorem. The difficult part is to show that $6 \rightarrow 1$. The proof is based on two main claims:

- If $\text{VCdim}(\mathcal{H}) = d$, then even though \mathcal{H} might be infinite, when restricting it to a finite set $C \subset \mathcal{X}$, its "effective" size, $|\mathcal{H}_C|$, is only $O(|C|^d)$. That is, the size of \mathcal{H}_C grows polynomially rather than exponentially with $|C|$. This claim is often referred to as *Shelah's lemma*, but it has also been stated and proved independently by Shelah and by Perles. The formal statement is given in Section 6.5.1 later.
- In Section 4 we have shown that finite hypothesis classes enjoy the uniform convergence property. In Section 6.5.2 later we generalize this result and show that uniform convergence holds whenever the hypothesis class has a "small effective size." By "small effective size" we mean classes for which $|\mathcal{H}_C|$ grows polynomially with $|C|$.

① → ③ Assume $\inf_{h \in \mathcal{H}} |\ell_S(h) - \ell(h)| \leq \epsilon, \forall h \in \mathcal{H} \Rightarrow 1 - \delta \text{ for } |S| \geq m_H(\epsilon, \delta) \quad h_S \text{ is the ERM rule}$

$$\ell_D(h_S) \leq \ell_S(h_S) + \epsilon \leq \ell_S(h) + \epsilon \leq \ell_D(h) + 2\epsilon.$$

$$\Rightarrow \ell_D(h_S) \leq \min_{h \in \mathcal{H}} \ell_D(h) + 2\epsilon$$

② → ③ Obvious

③ → ④ Realizability Obvious

② → ⑤ Obvious

② → ⑥

⑤ → ①

THEOREM 5.1 (No-Free-Lunch). Let A be any learning algorithm for the task of binary classification with respect to the 0-1 loss over a domain \mathcal{X} . Let m be any number smaller than $|\mathcal{X}|/2$, representing a training set size. Then, there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ such that:

- There exists a function $f: \mathcal{X} \rightarrow \{0, 1\}$ with $\ell_P(f) = 0$.
- With probability of at least $1/7$ over the choice of $S \sim \mathcal{D}^m$ we have that $\ell_D(A(S)) \geq 1/8$.