

# Machine Learning 6: VC dimension, Sauer's Lemma, Fundamental Theorem of Statistical Learning

Master 2 Computer Science

---

Aurélien Garivier

2019-2020



# Table of contents

1. VC dimension and Sauer's lemma
2. Finite VC dimension implies Uniform Convergence
3. Finite VC-dimension implies learnability

## **VC dimension and Sauer's lemma**

---

# Shattering

## Definition

Let  $\mathcal{H}$  be a class of functions  $\mathcal{X} \rightarrow \{0, 1\}$  and let  $C = \{c_1, \dots, c_m\} \subset \mathcal{X}$ . The *restriction* of  $\mathcal{H}$  to  $C$  is the set of functions  $C \rightarrow \{0, 1\}$  that can be derived from  $\mathcal{H}$ :

$$\mathcal{H}_C = \left\{ (c_1, \dots, c_m) \rightarrow (h(c_1), \dots, h(c_m)) : h \in \mathcal{H} \right\}.$$

## Shattering

A hypothesis class  $\mathcal{H}$  *shatters* a finite set  $C \subset \mathcal{X}$  if  $\mathcal{H}_C = \{0, 1\}^C$ .

Example:

- $\mathcal{H} = \{ \mathbb{1}_{(-\infty, a]} : a \in \mathbb{R} \}$ .
- $\mathcal{H}_{\text{rec}}^2 = \{ \mathbb{1}_{[a_1, b_1] \times [a_2, b_2]} : a_1 \leq b_1 \text{ and } a_2 \leq b_2 \}$ .

## Definition

The *Vapnik Chervonenkis dimension*  $\text{VCdim}(\mathcal{H})$  of a hypothesis class  $\mathcal{H}$  is the maximal size of a set  $C \subset \mathcal{X}$  that can be shattered by  $\mathcal{H}$ . If  $\mathcal{H}$  can shatter sets of arbitrarily large size we say that  $\text{VCdim}(\mathcal{H}) = \infty$ .

## Theorem

Let  $\mathcal{H}$  be a class of infinite VC-dimension. Then  $\mathcal{H}$  is not PAC-learnable.

**Proof:** for every training size  $m$ , there exists a set  $C$  of size  $2m$  that is shattered by  $\mathcal{H}$ . By the NFL theorem, for every learning algorithm  $A$  there exists a probability distribution  $D$  over  $\mathcal{X} \times \{0, 1\}$  such that  $L_D(h) = 0$  but with probability at least  $1/7$  over the training set, we have  $L_D(A(S)) \geq 1/8$ .

# Fundamental theorem of PAC learning

Let  $\mathcal{H}$  be a hypothesis class of functions from a domain  $\mathcal{X}$  to  $\{0, 1\}$  and let the loss function of 0 – 1 loss. Then the following propositions are equivalent:

1.  $\mathcal{H}$  has the uniform convergence property,
2. any ERM rule is a successful agnostic PAC learner for  $\mathcal{H}$ ,
3.  $\mathcal{H}$  is agnostic PAC learnable,
4.  $\mathcal{H}$  is PAC learnable,
5. any ERM rule is a successful PAC learner for  $\mathcal{H}$ ,
6.  $\mathcal{H}$  has finite VC-dimension.

# Fundamental theorem of PAC learning (quantitative version)

Let  $\mathcal{H}$  be a hypothesis class of functions from a domain  $\mathcal{X}$  to  $\{0, 1\}$  and let the loss function of 0 – 1 loss. Assume that  $\text{VCdim}(\mathcal{H}) < \infty$ . Then there exist constants  $C_1, C_2$  such that:

1.  $\mathcal{H}$  has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2},$$

2.  $\mathcal{H}$  is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2},$$

3.  $\mathcal{H}$  is PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}.$$

# Sauer's lemma

## Definition

Let  $\mathcal{H}$  be a hypothesis class. Then the *growth function* of  $\mathcal{H}$ , denoted  $\tau_{\mathcal{H}} : \mathbb{N} \rightarrow \mathbb{N}$ , is defined as the maximal number of different functions that can be obtained by restricting  $\mathcal{H}$  to a set of size  $m$ :

$$\tau_{\mathcal{H}}(m) = \max_{C \subset X: |C|=m} |\mathcal{H}_C|.$$

Note: if  $\text{VCdim}(\mathcal{H}) = d$ , then for any  $m \leq d$  we have  $\tau_{\mathcal{H}}(m) = 2^m$ .

## Sauer's lemma

Let  $\mathcal{H}$  be a hypothesis class with  $d = \text{VCdim}(\mathcal{H}) < \infty$ . Then, for all  $m \geq d$ ,

$$\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d.$$

Think of example:  $\mathcal{H} = \{\mathbb{1}_{(-\infty, a]} : a \in \mathbb{R}\}$  with  $d = \text{VCdim}(\mathcal{H}) = 1$ .



# Proof of Sauer's lemma 1/2

In fact we prove the stronger claim:

$$|\mathcal{H}_C| \leq |\{B \subset C : \mathcal{H} \text{ shatters } B\}| \leq \sum_{i=0}^d \binom{m}{i}.$$

where the last inequality holds since no set of size larger than  $d$  is shattered by  $\mathcal{H}$ . The proof is by induction.

**m=1:** The empty set is always considered to be shattered by  $\mathcal{H}$ . Hence, either  $|\mathcal{H}(C)| = 1$  and  $d = 0$ , inequality  $1 \leq 1$ , or  $d \geq 1$  and the inequality is  $2 \leq 2$ .

**Induction:** Let  $C = \{c_1, \dots, c_m\}$ , and let  $C' = \{c_2, \dots, c_m\}$ . We note functions like vectors, and we define

$$Y_0 = \left\{ (y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in \mathcal{H}_C \text{ or } (1, y_2, \dots, y_m) \in \mathcal{H}_C \right\}, \text{ and}$$
$$Y_1 = \left\{ (y_2, \dots, y_m) : (0, y_2, \dots, y_m) \in \mathcal{H}_C \text{ and } (1, y_2, \dots, y_m) \in \mathcal{H}_C \right\}.$$

Then  $|\mathcal{H}_C| = |Y_0| + |Y_1|$ . Moreover,  $Y_0 = \mathcal{H}_{C'}$  and hence by the induction hypothesis:

$$|Y_0| \leq |\mathcal{H}_{C'}| \leq |\{B \subset C' : \mathcal{H} \text{ shatters } B\}| = |\{B \subset C : c_1 \notin B \text{ and } \mathcal{H} \text{ shatters } B\}|$$

Next, define

$$\mathcal{H}' = \left\{ h \in \mathcal{H} : \exists h' \in \mathcal{H} \text{ s.t. } \forall 1 \leq i \leq n, h'(c_i) = \begin{cases} 1 - h(c_1) & \text{if } i = 1 \\ h(c_i) & \text{otherwise} \end{cases} \right\}$$

Note that  $\mathcal{H}'$  shatters  $B \subset C'$  iff  $\mathcal{H}$  shatters  $B \cup \{c_1\}$ , and that  $Y_1 = \mathcal{H}'_{C'}$ . Hence, by the induction hypothesis,

$$|Y_1| = |\mathcal{H}'_{C'}| \leq |\{B \subset C' : \mathcal{H}' \text{ shatters } B\}| = |\{B \subset C' : \mathcal{H}' \text{ shatters } B \cup \{c_1\}\}|$$
$$= |\{B \subset C : c_1 \in B \text{ and } \mathcal{H}' \text{ shatters } B\}| \leq |\{B \subset C : c_1 \in B \text{ and } \mathcal{H} \text{ shatters } B\}|.$$

Overall,

$$|\mathcal{H}_C| = |Y_0| + |Y_1| \leq |\{B \subset C : c_1 \notin B \text{ and } \mathcal{H} \text{ shatters } B\}| + |\{B \subset C : c_1 \in B \text{ and } \mathcal{H} \text{ shatters } B\}| = |\{B \subset C : \mathcal{H} \text{ shatters } B\}|.$$

## Proof of Sauer's lemma 2/2

For the last inequality, one may observe that if  $m \geq 2d$ , defining  $N \sim \mathcal{B}(m, 1/2)$ , Chernoff's inequality and inequality  $\log(u) \geq (u - 1)/u$  yield

$$\begin{aligned} -\log \mathbb{P}(N \leq d) &\geq m \operatorname{kl} \left( \frac{d}{m}, \frac{1}{2} \right) \geq d \log \frac{2d}{m} + (m - d) \log \frac{2(m - d)}{m} \\ &\geq m \log(2) + d \log \frac{d}{m} + (m - d) \frac{-d/m}{(m - d)/m} \\ &= m \log(2) + d \log \frac{d}{em}, \end{aligned}$$

and hence

$$\sum_{i=0}^d \binom{m}{i} = 2^m \mathbb{P}(N \leq d) \leq \exp \left( -d \log \frac{d}{em} \right) = \left( \frac{em}{d} \right)^d.$$

Besides, for the case  $d \leq m \leq 2d$ , the inequality is obvious since  $(em/d)^d \geq 2^m$ : indeed, function  $f : x \mapsto -x \log(x/e)$  is increasing on  $[0, 1]$ , and hence for all  $d \leq m \leq 2d$ :

$$\frac{d}{m} \log \frac{em}{d} = f(d/m) \geq f(1/2) = \frac{1}{2} \log(2e) \geq \log(2),$$

which implies

$$\left( \frac{em}{d} \right)^d = \exp \left( d \log \frac{em}{d} \right) \geq \exp(m \log(2)) = 2^m.$$

Alternately, you may simply observe that for all  $m \geq d$ ,

$$\left( \frac{d}{m} \right)^d \sum_{i=0}^d \binom{m}{i} \leq \sum_{i=0}^d \left( \frac{d}{m} \right)^i \binom{m}{i} \leq \sum_{i=0}^m \left( \frac{d}{m} \right)^i \binom{m}{i} = \left( 1 + \frac{d}{m} \right)^m \leq e^d.$$

# **Finite VC dimension implies Uniform Convergence**

---

# Finite VC dimension implies Uniform Convergence

## Theorem

Let  $\mathcal{H}$  be a class and let  $\tau_{\mathcal{H}}$  be its growth function. Then, for every distribution  $D$  and for every  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over the choice of the sample  $S \sim D^{\otimes m}$  we have

$$\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \leq \frac{1 + \sqrt{\log(\tau_{\mathcal{H}}(2m))}}{\delta \sqrt{m/2}}.$$

Note: this result is sufficient to prove that finite VC-dim  $\implies$  learnable, but the dependency in  $\delta$  is not correct at all: roughly speaking, the factor  $1/\delta$  can be replaced by  $\log(1/\delta)$ .

# Proof: symmetrization and Rademacher complexity (1/2)

We consider the 0-1 loss  $\ell(h, (x, y)) = \mathbb{1}\{h(x) \neq y\}$ , or any  $[0, 1]$ -valued loss  $\ell$ . We denote  $Z_i = (X_i, Y_i)$ , and observe that  $L_D(h) = \mathbb{E}_{Z_i}[\ell(h, Z_i)] = \mathbb{E}_{S'}[L_{S'}(h)]$  if  $S' = Z'_1, \dots, Z'_m$  denotes another iid sample of  $D$ . Hence,

$$\begin{aligned}\mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \right] &= \mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} |\mathbb{E}_{S'}[L_{S'}(h)] - L_S(h)| \right] = \mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{S'} [L_{S'}(h) - L_S(h)] \right| \right] \\ &\leq \mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} \mathbb{E}_{S'} \left[ |L_{S'}(h) - L_S(h)| \right] \right] \leq \mathbb{E}_S \left[ \mathbb{E}_{S'} \left[ \sup_{h \in \mathcal{H}} |L_{S'}(h) - L_S(h)| \right] \right] \\ &= \mathbb{E}_{S, S'} \left[ \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \ell(h, Z'_i) - \ell(h, Z_i) \right| \right] \\ &= \mathbb{E}_{S, S'} \left[ \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i (\ell(h, Z'_i) - \ell(h, Z_i)) \right| \right] \quad \text{for all } \sigma \in \{\pm 1\}^m \\ &= \mathbb{E}_\Sigma \mathbb{E}_{S, S'} \left[ \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \Sigma_i (\ell(h, Z'_i) - \ell(h, Z_i)) \right| \right] \quad \text{if } \Sigma \sim \mathcal{U}(\{\pm 1\}^m) \\ &= \mathbb{E}_{S, S'} \mathbb{E}_\Sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \Sigma_i (\ell(h, Z'_i) - \ell(h, Z_i)) \right| \right].\end{aligned}$$

Now, for every  $S, S'$ , let  $C = C_{S, S'}$  be the instances appearing in  $S$  and  $S'$ . Then  $\forall \sigma \in \{-1, 1\}^m$ ,

$$\sup_{h \in \mathcal{H}} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i (\ell(h, Z'_i) - \ell(h, Z_i)) \right| = \max_{h \in \mathcal{H}_C} \frac{1}{m} \left| \sum_{i=1}^m \sigma_i (\ell(h, Z'_i) - \ell(h, Z_i)) \right|.$$

## Proof: symmetrization and Rademacher complexity (2/2)

Moreover, for every  $h \in \mathcal{H}_C$  let  $Z_h = \frac{1}{m} \sum_{i=1}^m \Sigma_i(\ell(h, Z'_i) - \ell(h, Z_i))$ . Then  $\mathbb{E}_\Sigma[Z_h] = 0$ , each summand belongs to  $[-1, 1]$  and by Hoeffding's inequality, for every  $\epsilon > 0$ :

$$\mathbb{P}_\Sigma[|Z_h| \geq \epsilon] \leq 2 \exp\left(-\frac{m\epsilon^2}{2}\right).$$

Hence, by the union bound,

$$\mathbb{P}_\Sigma\left[\max_{h \in \mathcal{H}_C} |Z_h| \geq \epsilon\right] \leq 2|\mathcal{H}_C| \exp\left(-\frac{m\epsilon^2}{2}\right).$$

The following lemma permits to deduce that

$$\mathbb{E}_\Sigma\left[\max_{h \in \mathcal{H}_C} |Z_h|\right] \leq \frac{1 + \sqrt{\log(|\mathcal{H}_C|)}}{\sqrt{m/2}} \leq \frac{1 + \sqrt{\log(\tau_{\mathcal{H}}(2m))}}{\sqrt{m/2}}.$$

Hence,

$$\mathbb{E}_S\left[\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)|\right] \leq \mathbb{E}_{S, S'} \mathbb{E}_\Sigma\left[\sup_{h \in \mathcal{H}} \frac{1}{m} \left|\sum_{i=1}^m \Sigma_i(\ell(h, z'_i) - \ell(h, z_i))\right|\right] \leq \frac{1 + \sqrt{\log(\tau_{\mathcal{H}}(2m))}}{\sqrt{m/2}},$$

and we conclude by using Markov's inequality (poor idea! Better: McDiarmid's inequality).

# Technical Lemma

## Lemma

Let  $a > 0$ ,  $b > 1$ , and let  $Z$  be a real-valued random variable such that for all  $t \geq 0$ ,  $\mathbb{P}(Z \geq t) \leq 2b \exp\left(-\frac{t^2}{a^2}\right)$ . Then

$$\mathbb{E}[Z] \leq a \left( \sqrt{\log(b)} + \frac{1}{\sqrt{\log(b)}} \right).$$

**Proof:**

$$\begin{aligned} \mathbb{E}[Z] &\leq \int_0^\infty \mathbb{P}(Z \geq t) dt \leq a\sqrt{\log(b)} + \int_{a\sqrt{\log(b)}}^\infty 2b \exp\left(-\frac{t^2}{a^2}\right) \\ &\leq a\sqrt{\log(b)} + 2b \int_{a\sqrt{\log(b)}}^\infty \frac{t}{a\sqrt{\log(b)}} \exp\left(-\frac{t^2}{a^2}\right) \\ &= a\sqrt{\log(b)} + \frac{2b}{a\sqrt{\log(b)}} \times \frac{a^2}{2} \exp\left(-\frac{(a\sqrt{\log(b)})^2}{a^2}\right) \\ &= a\sqrt{\log(b)} + \frac{a}{\sqrt{\log(b)}}. \end{aligned}$$

NB: cutting at  $a\sqrt{\log(2b)}$  gives a better but less nice inequality for our use.

**Finite VC-dimension implies  
learnability**

---



## Application: Finite VC-dim classes are agnostically learnable

It suffices to prove that finite VC-dim implies the uniform convergence property. From Sauer's lemma, for all  $m \geq d/2$  we have  $\tau_{\mathcal{H}}(2m) \leq (2em/d)^d$ . With the previous theorem, this yields that with probability at least  $1 - \delta$ :

$$\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \leq \frac{1 + \sqrt{d \log(2em/d)}}{\delta \sqrt{m/2}} \leq \frac{1}{\delta} \sqrt{\frac{8d \log(2em/d)}{m}}$$

as soon as  $\sqrt{d \log(2em/d)} \geq 1$ . To ensure that this is at most  $\epsilon$ , one may choose

$$m \geq \frac{8d \log(m)}{(\delta\epsilon)^2} + \frac{8d \log(2e/d)}{(\delta\epsilon)^2}.$$

By the following lemma, it is sufficient that

$$m \geq \frac{32d \log\left(\frac{4d}{(\delta\epsilon)^2}\right)}{(\delta\epsilon)^2} + \frac{16d \log\left(\frac{2e}{d}\right)}{(\delta\epsilon)^2}.$$

# Technical Lemma

## Lemma

Let  $a > 0$ . Then

$$x \geq 2a \log(a) \implies x \geq a \log(x).$$

**Proof:** For  $a \leq e$ , true for every  $x > 0$ . Otherwise, for  $a \geq \sqrt{e}$  we have  $2a \log(a) \geq a$  and thus for every  $t \geq 2a \log(a)$ , as  $f : t \mapsto t - a \log(t)$  is increasing on  $[a, \infty)$ ,  $f(t) \geq f(2a \log(a)) = a \log(a) - a \log(2 \log(a)) \geq 0$ , since for every  $a > 0$  it holds that  $a \geq 2 \log(a)$ .

## Lemma

Let  $a \geq 1, b > 0$ . Then

$$x \geq 4a \log(2a) + 2b \implies x \geq a \log(x) + b.$$

**Proof:** It suffices to check that  $x \geq 2a \log(x)$  (given by the above lemma) and that  $x \geq 2b$  (obvious since  $4a \log(2a) \geq 0$ ).