

II / Application to classes \mathcal{H} with finite VC-dimension

Recap: we proved the following result

Theorem: the ERM $\hat{h}_n \in \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq h(X_i)}$ satisfies

$$L(\hat{h}_n) - \inf_{h \in \mathcal{H}} L(h) \leq 4 \sqrt{\frac{2 \lg(25 S_n(\mathcal{H}))}{n}} + \sqrt{\frac{2 \lg 1/\delta}{n}} \quad \text{with prob} \geq 1 - \delta \quad (1)$$

$$\mathbb{E}[L(\hat{h}_n)] - \inf_{h \in \mathcal{H}} L(h) \leq 4 \sqrt{\frac{2 \lg(25 S_n(\mathcal{H}))}{n}} \quad (2)$$

where $L(h) = \mathbb{P}(Y \neq h(X)) = \int \mathbb{1}_{y \neq h(x)} d\mathbb{P}^{(X,Y)}$ ("risk of h ")

and $S_n(\mathcal{H}) = \max_{(x_1, \dots, x_n) \in \mathcal{X}^n} |\{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}|$ ("shattering-coefficient")

We proved directly (1) but (2) is even more direct.

Next we show a special case where the shattering-coefficient $S_n(\mathcal{H})$ can be upper bounded in a simple fashion. We need a new notion:

Definition: the Vapnik-Chervonenkis dimension of \mathcal{H} (or "VC-dimension") is defined by

$$V := \sup \{k \in \mathbb{N} : S_k(\mathcal{H}) = 2^k\} \quad (S_0(\mathcal{H}) := 1 \text{ by convention})$$
$$= \sup \{k \in \mathbb{N} : \exists (x_1, \dots, x_k) \in \mathcal{X}^k, |\{(h(x_1), \dots, h(x_k)) : h \in \mathcal{H}\}| = 2^k\}$$

$$\in \mathbb{N} \cup \{+\infty\}.$$

It is the maximal number of points that can be labelled in all possible ways by classifiers in \mathcal{H} .

- Examples:
- affine classifiers: $\mathcal{H} = \{h(x) = \text{sign}(\langle a, x \rangle + b), (a, b) \in \mathbb{R}^d \times \mathbb{R}\}$
 - convex polygonal classifiers: $\mathcal{H} = \{h(x) = 2\mathbb{1}_A(x) - 1 : A \text{ convex polytope of } \mathbb{R}^d\}$
- $\chi = \mathbb{R}^d$
in both examples
- has VC-dimension equal to $d+1$
 - has VC-dimension equal to $+\infty$.

A useful combinatorial result is:

Lauer's lemma

Let $\mathcal{H} \subseteq \{-1, 1\}^{\mathcal{X}}$ be a set of classifiers with finite VC-dimension $V \in \mathbb{N}$.

Then, the shattering coefficients $S_m(\mathcal{H})$ can be upper bounded as:

$$S_m(\mathcal{H}) \leq \sum_{k=0}^V \binom{m}{k} \leq (m+1)^V$$

↑ a tighter upper bound of $(\frac{em}{V})^V$ can also be proved.

Substituting the bound in Lauer's lemma into the excess risk bounds (1) and (2), we obtain:

Corollary: if \mathcal{H} has VC-dimension $V \in \mathbb{N}^*$, then the ERM \hat{h}_m satisfies

$$L(\hat{h}_m) - \inf_{h \in \mathcal{H}} L(h) \leq 8 \sqrt{\frac{V \lg(m+1)}{n}} + \sqrt{\frac{2 \lg(1/\delta)}{n}} \text{ with prob } \geq 1-\delta \quad (1')$$

$$\mathbb{E}[L(\hat{h}_m)] - \inf_{h \in \mathcal{H}} L(h) \leq 8 \sqrt{\frac{V \lg(m+1)}{n}} \quad (2')$$

In the next section we address the optimality of (2').

III - Minimax lower bound for the misclassification excess risk with VC-classes

In this section, we prove that the risk bound of order $\sqrt{\frac{V \log m}{n}}$ is optimal in the worst case (up to log factors).

Theorem 3

We consider the classification problem with input space \mathcal{X} and output space $\{-1, 1\}$.

Let $\mathcal{G} \subseteq \{-1, 1\}^{\mathcal{X}}$ be a class of classifiers with finite VC-dimension $V \in \mathbb{N}^*$.

Then, for all $n \geq c_1 V$,

$$\inf_{\hat{g}_n} \sup_{P \in \mathcal{U}_1^+(\mathcal{X} \times \{-1, 1\})} \left\{ \mathbb{E}[R(\hat{g}_n)] - \inf_{g \in \mathcal{G}} R(g) \right\} \geq c_2 \sqrt{\frac{V}{n}}$$

where $c_1, c_2 > 0$ are absolute constants.

Our proof is inspired from Bassat and Kédelec (AoS'06) who actually proved a finer lower bound in terms of the margin $h = \text{ess inf} |\mathbb{P}(Y=1|X) - 1/2|$. Here, we don't care about the margin, but the interested reader can note that we consider the particular case $h \approx \sqrt{\frac{V}{n}}$.

Proof: As in Section II, we restrict the sup over an appropriately chosen finite set, which we define next:

- Since \mathcal{G} has VC-dimension equal to V , there exists $\{x_1, \dots, x_V\} \subseteq \mathcal{X}$ such that \mathcal{G} shatters $\{x_1, \dots, x_V\}$, i.e., $|\{(g(x_1), \dots, g(x_V)) : g \in \mathcal{G}\}| = 2^V$. In particular, $x_i \neq x_j$ for all $1 \leq i \neq j \leq V$.
- Denote by μ the uniform probability measure over $\{x_1, \dots, x_V\}$.

- Define, for all $b \in \{-1, 1\}^V$, the function P_b on $\{x_1, \dots, x_V\}$ by

$$P_b(x_i) = \frac{1}{2}(1 + b_i h) \quad , \quad 1 \leq i \leq V.$$

where $h \in (0, 1)$ will be chosen later.

Now, for all $b \in \{-1, 1\}^V$, define P_b as the joint probability distribution on $X \times \{-1, 1\}$ such that, under P_b , X has distribution μ and $Y|X=x \sim P_b(x) \delta_1 + (1 - P_b(x)) \delta_{-1}$, for all $x \in \{x_1, \dots, x_V\}$.

(In other words: for all nonnegative measurable mapping $\varphi: X \rightarrow \mathbb{R}$,

$$\mathbb{E}_P[\varphi(X) \mathbb{1}_{\{Y=1\}}] = \frac{1}{V} \sum_{i=1}^V \varphi(x_i) P_b(x_i)$$

$$\mathbb{E}_P[\varphi(X) \mathbb{1}_{\{Y=-1\}}] = \frac{1}{V} \sum_{i=1}^V \varphi(x_i) (1 - P_b(x_i))$$

Note that the Bayes classifier $t_b^*: X \rightarrow \{-1, 1\}$ under P_b is given by

$$\begin{cases} +1 & \text{if } P_b(x_i) \geq \frac{1}{2} \\ -1 & \text{if } P_b(x_i) < \frac{1}{2} \end{cases} \quad , \quad 1 \leq i \leq V$$

i.e. $t_b^*(x_i) = b_i$ for all $1 \leq i \leq V$.

Note also that, since \mathcal{G} shatters $\{x_1, \dots, x_V\}$:

$$\forall b \in \{-1, 1\}^V, \exists g_b \in \mathcal{G} : (g_b(x_1), \dots, g_b(x_V)) = b.$$

Therefore, g_b is the Bayes classifier under P_b , so that

$$R(g_b) = \inf_{\substack{\text{measurable} \\ g: X \rightarrow \{-1, 1\}}} R(g) = \inf_{g \in \mathcal{G}} R(g) \quad (5)$$

We are now ready to derive the lower bound.

Let $\Theta = 2^{\Gamma-1} = \{2x-1 : x \in \Gamma\} \subseteq \{-1, 1\}^V$ be the set provided by Vapnik-Chervakov-Gilbert's lemma rescaled to $\{-1, 1\}^V$.

We assume that $V \geq 6$ so that $|\Theta| \geq e^{6/8} \geq 2$.

$$\inf_{\hat{g}_m} \sup_{P \in \mathcal{P}_\beta^+(\mathcal{X} \times \{-1,1\})} \left\{ \mathbb{E}[R(\hat{g}_m)] - \inf_{g \in \mathcal{G}} R(g) \right\}$$

$$\geq \inf_{\hat{g}_m} \max_{b \in \mathcal{H}} \left\{ \mathbb{E}[R(\hat{g}_m)] - \inf_{g \in \mathcal{G}} R(g) \right\}$$

$$= \inf_{\substack{\text{measurable} \\ g: \mathcal{X} \rightarrow \{-1,1\}}} R(g) \text{ by (5)}$$

denotes expectation when P_b is the underlying prob. distribution

$$= \mathbb{E}_b \left[|2P_b(x) - 1| \cdot \mathbb{1}_{\{\hat{g}_m(x) \neq g_b(x)\}} \right]$$

$$\geq h \inf_{\hat{g}_m} \max_{b \in \mathcal{H}} \mathbb{E}_b \left[\frac{|\hat{g}_m(x) - g_b(x)|}{2} \right] \quad \text{since } |2P_b - 1| \geq h \text{ by construction}$$

$$\text{and } \mathbb{1}_{u \neq v} = \frac{|u-v|}{2} \text{ for all } u, v \in \{-1,1\}$$

$$= \frac{h}{2} \inf_{\hat{g}_m} \max_{b \in \mathcal{H}} \mathbb{E}_b \left[\|\hat{g}_m - g_b\|_{L^1(\mu)} \right] \quad (6)$$

Let $\hat{b}_m \in \operatorname{argmin}_{b \in \mathcal{H}} \|\hat{g}_m - g_b\|_{L^1(\mu)}$ (\hat{b}_m is a random variable)

As previously, we get from the triangle inequality that

$$\{\hat{b}_m \neq b\} \subseteq \left\{ \|\hat{g}_m - g_b\|_{L^1(\mu)} \geq \frac{\varepsilon}{2} \right\} \quad \text{where } \varepsilon > 0 \text{ is such that } \|g_{b_1} - g_{b_2}\|_{L^1(\mu)} > \varepsilon \text{ for all } b_1 \neq b_2 \in \mathcal{H}.$$

$$\begin{aligned} \text{(Otherwise, we would have } \|\hat{g}_m - g_{\hat{b}_m}\|_{L^1(\mu)} &\leq \|\hat{g}_m - g_{\hat{b}_m}\|_{L^1(\mu)} + \|\hat{g}_m - g_b\|_{L^1(\mu)} \\ &\leq 2 \|\hat{g}_m - g_{\hat{b}_m}\|_{L^1(\mu)} \text{ by def of } \hat{b}_m \\ &< \varepsilon, \text{ which is impossible.)} \end{aligned}$$

Here, we can take $\varepsilon = \frac{1}{2}$ because

[see below]

$\forall b \neq b' \in \Theta$, $b = 2x - 1$, $b' = 2x' - 1$, $x \neq x' \in \Gamma$,

$$\begin{aligned} \|g_b - g_{b'}\|_{L^1(\mu)} &= \frac{1}{V} \sum_{i=1}^V |g_b(x_i) - g_{b'}(x_i)| \\ &= \frac{1}{V} \sum_{i=1}^V |b_i - b'_i| \quad \text{by def of } g_b \text{ and } g_{b'} \\ &= \frac{c}{V} \sum_{i=1}^V \mathbb{1}_{\{b_i \neq b'_i\}} \\ &> \frac{c}{V} \cdot \frac{V}{4} = \frac{1}{c} \end{aligned}$$

Getting back to (6), we obtain

$$\begin{aligned} \inf_{\hat{g}_m} \sup_P \{E[R(\hat{g}_m)] - \inf_{g \in \mathcal{G}} R(g)\} &\geq \frac{h}{2} \inf_{\hat{g}_m} \max_{b \in \Theta} \frac{\varepsilon}{c} \mathbb{P}_b(\hat{b}_m \neq b) \\ &= \frac{h}{8} \left(1 - \sup_{\hat{g}_m} \min_{b \in \Theta} \mathbb{P}_b(\hat{b}_m = b)\right) \end{aligned}$$

Fix \hat{g}_m . We use Fano's inequality to upper bound $\min_{b \in \Theta} \mathbb{P}_b(\hat{b}_m = b)$:

- We have $|\Theta| \geq 2$ since we assumed that $V \geq 6$.
- Let $\tilde{b} \in \Theta$. We have: $\forall b \in \Theta \setminus \{\tilde{b}\}$,

$$KL(P_b^{\otimes m}, P_{\tilde{b}}^{\otimes m}) = m \cdot KL(P_b, P_{\tilde{b}})$$

$$= m \left(KL(\mu, \mu) + \int_{\mathcal{X}} d\mu(x) KL(\mathcal{B}(P_b(x)), \mathcal{B}(P_{\tilde{b}}(x))) \right)$$

by the chain rule for the Kullback-Leibler entropy,
and since under P_b , $Y|X=x \sim P_b(x) \delta_1 + (1-P_b(x)) \delta_{-1}$

$$= \frac{m}{V} \sum_{i=1}^V \underbrace{KL(\mathcal{B}(P_b(x_i)), \mathcal{B}(P_{\tilde{b}}(x_i)))}_{\substack{\text{if } b_i \neq \tilde{b}_i, \text{ this quantity} \\ \text{is equal to either} \\ KL(\mathcal{B}(\frac{1+b}{2}), \mathcal{B}(\frac{1-\tilde{b}}{2})) \text{ or} \\ KL(\mathcal{B}(\frac{1-\tilde{b}}{2}), \mathcal{B}(\frac{1+b}{2}))}} \mathbb{1}_{\{b_i \neq \tilde{b}_i\}}$$

$$\begin{aligned} &\uparrow P_b(x_i) \neq P_{\tilde{b}}(x_i) \\ &\Leftrightarrow b_i \neq \tilde{b}_i \end{aligned}$$

i.e. either $KL(B(q), B(1-q))$ or $KL(B(1-q), B(q))$, with $q := \frac{1+h}{2}$

But these two quantities are equal to $(2q-1) \lg \frac{q}{1-q} = h \lg \left(\frac{1+h}{1-h} \right)$

Hence:

$$\begin{aligned} KL(P_{\vec{b}}^{\otimes n}, P_{\vec{\tilde{b}}}^{\otimes n}) &= \frac{n}{V} \sum_{i=1}^V h \lg \left(\frac{1+h}{1-h} \right) \mathbb{1}_{\{b_i \neq \tilde{b}_i\}} \\ &\leq \frac{2nh^2}{1-h} \cdot \frac{1}{V} \sum_{i=1}^V \mathbb{1}_{\{b_i \neq \tilde{b}_i\}} \\ &\leq 4nh^2 \quad \text{if } h \in (0, \frac{1}{2}]. \quad (\text{condition to be checked later}) \end{aligned}$$

Applying Fano's inequality, we get: [see footnote]

$$\begin{aligned} \min_{b \in \mathbb{H}} P_b(\hat{b}_n = b) &\leq \max \left\{ \frac{2e}{2e+1}, \frac{4nh^2}{\ln |\mathbb{H}|} \right\} \\ &\leq \max \left\{ \frac{2e}{2e+1}, \frac{32nh^2}{V} \right\} \quad \text{since } \ln |\mathbb{H}| \geq \frac{V}{8} \end{aligned}$$

Choosing $h = \frac{1}{4} \sqrt{\frac{e}{2e+1} \cdot \frac{V}{n}}$, we get $h \in (0, \frac{1}{2}]$ when $n \geq \frac{e}{4(2e+1)} \cdot V$, and:

$$\begin{aligned} \inf_{\vec{g}} \sup_{P \in \mathcal{M}_1^+(\mathcal{X}^{\times V})} \{ E[R(\hat{g}_n)] - \inf_{g \in \mathcal{G}} R(g) \} &\geq \frac{h}{8} \left(1 - \frac{2e}{2e+1} \right) \\ &\geq \frac{1}{32} \sqrt{\frac{2e}{(2e+1)^3} \cdot \frac{V}{n}} \end{aligned}$$

The case $V \in \{1, 2, \dots, 5\}$ can be handled similarly (easier actually!) via Pinsker's inequality. The constants c_1 and c_2 may change a little. \square

Remark: I used a slightly different version of Fano's inequality than the one we studied in Lecture 2, namely:

Fano's inequality (Bingé's version)

Assume the events A_1, \dots, A_N form a partition of Ω (with $N \geq 2$)

• P_1, \dots, P_N are probability distributions on (Ω, \mathcal{F}) .

Then,

$$\min_{1 \leq j \leq N} P_j(A_j) \leq \max_{\vec{g}} \left\{ \frac{2e}{2e+1}, \frac{\overline{KL}}{\lg N} \right\} \quad \text{where } \overline{KL} := \frac{1}{N-1} \sum_{j=2}^N KL(P_j, P_1).$$

Exercise: prove a similar lower bound with the version of Fano involving $h_2(2)$.