# 17. Fundamental Thm of ML

(for binary classification)

**Sauer Lemma** ( aka Sauer - Shelah - Perles) (Lemma 6.10 [SS] or Thm 3.17 / Cor. 3.18 Mohri )

let $d = VCdim(\mathcal{H})$, then $\forall m \in \mathbb{N}$, $\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^{d} \binom{m}{i}$  $\left( \text{define } \binom{k}{\lambda} = 0 \text{ if } k < \lambda \right)$

Furthermore, if $m \leq d$ (or $d = \infty$) this bound is vacuous

(ie. it's $2^m$)

but if $\underline{m > d}$,

$e = $ Euler's constant

$$\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^{d} \binom{m}{i} \leq \left( \frac{e \, m}{d} \right)^{d}, \text{ a } \underline{\text{polynomial}} \text{ in } m \text{ (vs. } 2^m \text{)}$$

## proof sketch

note binomial thm, $(1+x)^m = \sum_{i=0}^{m} \binom{m}{i} x^i$   so $(x=1)$  $2^m = \sum_{i=0}^{m} \binom{m}{i}$

Mohri (Cor. 3.18) uses binomial thm plus $(1-x) \leq e^{-x}$

[SS] uses Stirling's approx. for $n!$ and induction, see Lemma A.5

& Lemma 6.10

## Relating back ...

**Thm 3.3 Mohri** / Thm 26.5 [SS], $Y \leq [0,1]$  $\forall \delta > 0$, w.p. $\geq 1-\delta$ (over $m$ iid samples in $S$)

$$\forall f \in \mathcal{F}, \quad \mathbb{E}_{z \sim D} f(z) - \frac{1}{m} \sum_{i=1}^{m} f(z_i) \leq \begin{cases} 2 \, \mathcal{R}_m(\mathcal{F}) + \sqrt{\log(\delta^{-1})/2m} \\ 2 \, \hat{\mathcal{R}}_s(\mathcal{F}) + 3\sqrt{\log(2\delta^{-1})/2m} \end{cases}$$

based on $\text{Rep}_D(\mathcal{F}, S) = \sup_{f \in \mathcal{F}} \mathbb{E} f(z) - \hat{\mathbb{E}}_s f(z)$

so w/ $\mathcal{F} = \ell \cdot \mathcal{H}$ for binary loss $\ell$

**Thm 3.5 Mohri** $Y = \{\pm 1\}$, $D$ any thing, $\forall \delta > 0$, w.p. $\geq 1-\delta$ over $S$ ($m$ iid samples)

$$\forall h \in \mathcal{H}, \quad L_D(h) \leq \hat{L}_s(h) + \begin{cases} \mathcal{R}_m(\mathcal{H}) + \sqrt{\log(\delta^{-1})/2m} \\ \hat{\mathcal{R}}_s(\mathcal{H}) + 3\sqrt{\log(2\delta^{-1})/2m} \end{cases}$$

and, implication of **Massart's Lemma**

$$\mathcal{R}_m(\mathcal{H}) \leq \sqrt{2 \log(\tau_{\mathcal{H}}(m))/m}$$

(or... use uniform convergence, Thm 6.11 [SS]:

$$\forall D, \forall \delta > 0, wp \geq 1-\delta, \forall h \in \mathcal{H} \quad | L_D(h) - \hat{L}_s(h) | \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2m))}}{\delta \cdot \sqrt{2m}} \quad )$$

So just need to bound **Rademacher Complexity** :

So combine **Massart** with **Sauer** to get      ( letting $d = VCdim(\mathcal{H})$ )

$$\mathcal{R}_m(H) \leq \sqrt{2 \log((\tfrac{em}{d})^d)/m} = \sqrt{\frac{2d \log(em/d)}{m}} \leq \sqrt{\frac{2 \, VCdim(H) \log(e \cdot m)}{m}}$$

( for binary classification )

§6.4 [SS]

**Thm 6.7** "Fundamental Thm. of Statistical (or PAC) learning" (for binary classification)   *Qualitative Version*

For $Y = \{0,1\}$ and the 0-1 loss function, the following are equivalent:

① $\mathcal{H}$ has the uniform convergence property

② Any ERM rule is a successful (agnostic) PAC learner

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

proof outline   ① → ② → ③ → ④ → ⑤ → ⑥

## Remarks

- for general learning (any loss function), uniform convergence ⟹ agnostic PAC learner
  For binary classif. (0-1 loss), vice-versa is true also!

- Some variants apply to regression ($l^1$ or $l^2$ loss) but not all learning tasks have such theorems

- See Thm 6.8 [SS] for a quantitative version
  i.e. agnostic PAC learnable w, $m_{\mathcal{H}}(\epsilon, \delta) \leq C \cdot \dfrac{VCdim(\mathcal{H}) + \log(1/\delta)}{\epsilon^2}$
  and this is tight up to a constant.

- For binary classif., $VCdim < \infty$ iff PAC learnable ... pretty neat!

## Proofs

-- mostly follow from our previous results

⑥ ⟹ ① via Massart's Lemma + Sauer lemma to bound Rademacher complexity, and this bounds representativeness which is basically what's needed for uniform convergence. See Thm 6.11 [SS], use Markov's Ineq. too

**History**  Vapnik + Chervonenkis '71
  As necessary condition for PAC, see Blumer, Andrzej Ehrenfeucht, David Haussler
  & Manfred Warmuth '89

*student of* →
*CU CS faculty (emeritus)*
*Founding member of CS dept.*