

FOSML: Lecture 7

①

The bound of the growth function in terms of VC dimension immediately yields a gen. bound in terms of VC dimension:

Corollary (VC dimension generalization bounds): Let \mathcal{H} be a family of functions taking values in $\{-1, 1\}$, with VC dimension d . Then $\forall \delta > 0$, with probability exceeding $1 - \delta$,

$$\forall h \in \mathcal{H}, R(h) \leq \hat{R}_h(S) + \sqrt{\frac{2d \log\left(\frac{em}{d}\right)}{m}} + \sqrt{\frac{\log(1/\delta)}{2m}}.$$

$$= \hat{R}_h(S) + O\left(\sqrt{\frac{\log(m/d)}{m/d}}\right).$$

In particular, as $d \rightarrow \infty$, our generalization bound blows up. This suggests a "principle of parsimony": all else being equal, we ought to prefer the simplest explanation, e.g. hypothesis class of smaller VC dimension.

A VC-analysis also yields lower bounds on the generalization error, which may be interpreted as a ~~measure~~ measure of fundamental hardness.

Recall that $R_D(h_S, f)$ is the generalization error for a function $f \in \mathcal{H}$ with respect to the predictor h_S for a sample $S \sim D$.

Theorem (lower bound, realizable case): Let \mathcal{H} be a hypothesis class with VC-dimension $d \geq 1$. Then $\forall m \geq 1$, and any learning algorithm A , there exists a distribution D over X and a target function $f \in \mathcal{H}$ such that $\mathbb{P}_{S \sim D} \left(R_D(h_S, f) > \frac{d-1}{32m} \right) \geq \frac{1}{100}$.

Proof: Since $V(\dim(H)=d, \exists$ a set $\bar{X} = \{x_0, x_1, \dots, x_{d-1}\}$ that is shattered by H .
 For any $\epsilon > 0$, we choose D such that its support is reduced to the shattered set \bar{X} and so that one point (WLOG, let it be x_0) is sampled with probability $1-8\epsilon$, and all other points (x_1, \dots, x_{d-1}) with probability $\frac{8\epsilon}{d-1}$. That is, D is s.t.

- (a.) $\mathbb{P}_{x \sim D}(x=x_0) = 1-8\epsilon,$
- (b.) $\mathbb{P}_{x \sim D}(x=x_i, i=1, \dots, d-1) = \frac{8\epsilon}{d-1}$

For this choice of D , most samples would contain x_0 , and A cannot do better than flip a coin on the points $\{x_i\}_{i=1}^{d-1}$ which do not appear in the training set. We now proceed to show that sufficiently many of these points $\{x_i\}_{i=1}^{d-1}$ are not in the sample, and thus that the generalization error is lower bounded in terms of d .

Suppose WLOG that A makes no error on x_0 . Let $S \sim D$ be a sample, and let $\bar{S} = S \cap \{x_1, \dots, x_{d-1}\}$. Let \mathcal{S} be the set of samples S of size n s.t. $|\bar{S}| \leq \frac{d-1}{2}$.
 Now, fix a sample $S \in \mathcal{S}$ and let ν be the uniform distribution over all labelings $f: \bar{X} \rightarrow \{0,1\}$, which are all in H since \bar{X} is shattered by H . We estimate:

$$\begin{aligned} \mathbb{E}_{f \sim \nu} (R_D(h_S, f)) &= \sum_f \cancel{P_D(x)} R_D(h_S, f) P_\nu(f) \\ &= \sum_f \sum_{x \in \bar{X}} \mathbb{1}_{\{h_S(x) \neq f(x)\}} P_D(x) P_\nu(f) \\ &\geq \sum_f \sum_{x \in \bar{S}} \mathbb{1}_{\{h_S(x) \neq f(x)\}} P_D(x) P_\nu(f) \\ &= \sum_{x \in \bar{S}} \left(\sum_f \mathbb{1}_{\{h_S(x) \neq f(x)\}} P_\nu(f) \right) P_D(x) \end{aligned}$$

$$= \sum_{x \in \bar{S}} \left(\frac{1}{2}\right) P_D(x)$$

Since all possible labelings for t are considered uniformly by ν

$$= \frac{1}{2} \sum_{x \in \bar{S}} P_D(x)$$

$$\geq \frac{1}{2} \cdot \underbrace{\left(\frac{d-1}{2}\right)}_{|\bar{S}|} \cdot \underbrace{\left(\frac{8\epsilon}{d-1}\right)}_{P_D(x=x_i)} = 2\epsilon$$

upper bound on $P_D(x=x_i)$

So, we have shown that $\forall S \in \mathcal{S}, \mathbb{E}_{f \sim \nu} (R_D(h_S, f)) \geq 2\epsilon$

$$\Rightarrow \mathbb{E}_{S \sim \mathcal{S}} \left(\mathbb{E}_{f \sim \nu} (R_D(h_S, f)) \right) \geq 2\epsilon$$

Fubini Theorem: $\Rightarrow \mathbb{E}_{f \sim \nu} \left(\mathbb{E}_{S \sim \mathcal{S}} (R_D(h_S, f)) \right) \geq 2\epsilon$
 $\int_S \int_X = \int_X \int_S$ under certain conditions

In particular, $\mathbb{E}_{S \sim \mathcal{S}} (R_D(h_S, f_0)) \geq 2\epsilon$ for some $f_0 \in \mathcal{H}$. Note that

$R_D(h_S, f_0) \leq P_D(\bar{X} - \{x_0\})$, so that we may analyze this expected value as

$$\mathbb{E}_{S \in \mathcal{S}} (R_D(h_S, f_0)) = \sum_{S: R_D(h_S, f_0) \geq \epsilon} R_D(h_S, f_0) P(R_D(h_S, f_0)) + \sum_{S: R_D(h_S, f_0) < \epsilon} \dots$$

$$\leq \mathbb{P}_D(\bar{X} \setminus \{x_0\}) \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon) + \varepsilon \sum_{S \in \mathcal{S}} \mathbb{P}[R_D(h_S, f_0) < \varepsilon]$$

$$\leq 8\varepsilon \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon) + \varepsilon (1 - \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon))$$

$$= 7\varepsilon \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon) + \varepsilon$$

We have thus shown that $2\varepsilon \leq 7\varepsilon \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon) + \varepsilon$

$$\Rightarrow \frac{1}{7} \leq \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon)$$

Thus, $\mathbb{P}_S(R_D(h_S, f_0) \geq \varepsilon) \geq \sum_{S \in \mathcal{S}} \mathbb{P}(R_D(h_S, f_0) \geq \varepsilon) \mathbb{P}(S \in \mathcal{S})$
 $\geq \frac{1}{7} \mathbb{P}(S \in \mathcal{S}).$

We will conclude by lowering bounding $\mathbb{P}(S) = \mathbb{P}(S \in \mathcal{S})$. Recall the multiplicative Chernoff bound: "for any X_1, \dots, X_m independent r.v. distributed according to D with mean μ and support $[0, 1]$, $\forall \gamma \in [0, \frac{1}{2} - 1]$, the following hold for $\hat{\rho} = \frac{1}{m} \sum_{i=1}^m X_i$:
 $\mathbb{P}(\hat{\rho} \geq (1+\gamma)\mu) \leq \exp(-m\mu\gamma^2/3)$; $\mathbb{P}(\hat{\rho} \leq (1-\gamma)\mu) \leq \exp(-m\mu\gamma^2/2)$." Applying this, we see: $1 - \mathbb{P}(S) = \mathbb{P}(S_m \geq 8\varepsilon(1+\gamma)) \leq \exp(-8\varepsilon m \gamma^2/3)$. Setting $\varepsilon = (d-1)/32m$ and $\gamma = 1$ yields

$$\mathbb{P}(S_m \geq \frac{d-1}{2}) \leq \exp(-(d-1)/12) \leq \exp(1/12) \leq 1-7\delta \text{ for } \delta \leq 1/100. \text{ We conclude } \mathbb{P}(S) \leq 7\delta \Rightarrow \mathbb{P}_S(R_D(h_S, f_0) \geq \varepsilon) \geq \delta, \text{ as desired.}$$