

Machine Learning Theory (CS 6783)

Lecture 5 : Symmetrization, Growth Function, and Effective Size

1 Recap

Last class we showed that

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq \sup_D \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right]$$

This was using the Empirical Risk Minimizer (ERM)

1. When $|\mathcal{F}| < \infty$, using the above we showed that

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq \sqrt{\frac{\log |\mathcal{F}|}{n}}$$

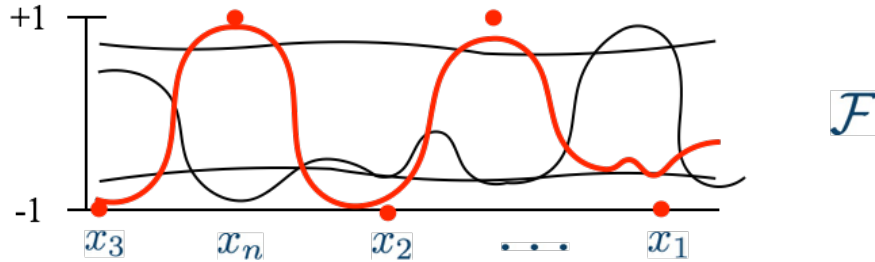
2. For countably infinite class we showed MDL bound and the algorithm based on this bound.
3. However the learning rate was not uniform over \mathcal{F}

2 Symmetrization and Rademacher Complexity

$$\begin{aligned} & \mathbb{E}_S [L_D(\hat{y}_{\text{erm}})] - \inf_{f \in \mathcal{F}} L_D(f) \\ & \leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ & \leq \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ & = \mathbb{E}_{S, S'} \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t (\ell(f(x'_t), y'_t) - \ell(f(x_t), y_t)) \right\} \right] \\ & \leq 2 \mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right\} \right] \\ & =: \mathcal{R}_n(\mathcal{F}) \end{aligned}$$

Where in the above each ϵ_t is a Rademacher random variable that is +1 with probability 1/2 and -1 with probability 1/2. The above is called Rademacher complexity of the loss class $\ell \circ \mathcal{F}$. In general Rademacher complexity of a function class measures how well the function class correlates with random signs. The more it can correlate with random signs the more complex the class is.

Example : $\mathcal{X} = [0, 1]$, $\mathcal{Y} = [-1, 1]$



3 Infinite \mathcal{F} : Binary Classes and Growth Function

First let us simplify the Rademacher complexity for binary classification problem. Note that for binary classification problem where $\mathcal{Y} \in \{\pm 1\}$, the loss can be rewritten as

$\ell(y', y) = \mathbf{1}_{\{y \neq y'\}} = \frac{1 - y \cdot y'}{2}$. Hence

$$\begin{aligned} 2\mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right\} \right] &= 2\mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t \frac{1 - f(x_t) \cdot y_t}{2} \right\} \right] \\ &= \mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t y_t f(x_t) \right] \end{aligned}$$

Now consider the inner term in the expectation above, ie. $\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t y_t f(x_t) \right]$. Note that given any fixed choice of $y_1, \dots, y_n \in \{\pm 1\}$, $\epsilon_1 y_1, \dots, \epsilon_n y_n$ are also Rademacher random variables. Hence for the binary classification problem,

$$2\mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right\} \right] = \mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t f(x_t) \right]$$

In the above statement we moved from Rademacher complexity of loss class $\ell \circ \mathcal{F}$ to the Rademacher complexity of the function class \mathcal{F} for binary classification task. This is a precursor to what we will refer to as contraction lemma which we will show later.

Why is symmetrization useful? Think what we gain for an infinite class ...

4 Effective size of function class on Data

Why is the introduction of Rademacher averages important ? To analyze the term, $\mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right]$ consider the inner expectation, that is conditioned on sample consider the term $\mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right]$. Note that $\frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t)$ is still average of 0 mean random variables (conditioned on data) and we can apply Hoeffding bound for each

fixed $f \in \mathcal{F}$ individually. Now \mathcal{F} might be an infinite class, but, conditioned on input instances $(x_1, y_1), \dots, (x_n, y_n)$, one can ask, what is the size of the projection set

$$\mathcal{F}_{|x_1, \dots, x_n} = \{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\}$$

For any binary class \mathcal{F} , first note that this set can have a maximum cardinality of 2^n however it could be much smaller. In fact we can have,

$$\mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right] = \mathbb{E}_S \mathbb{E}_\epsilon \left[\sup_{\mathbf{f} \in \mathcal{F}_{|x_1, \dots, x_n}} \frac{1}{n} \sum_{t=1}^n \epsilon_t \ell(\mathbf{f}[t], y_t) \right] \leq \mathbb{E}_S \left[\sqrt{\frac{\log |\mathcal{F}_{|x_1, \dots, x_n}|}{n}} \right]$$

where the last step is using the finite Lemma. Now one can define the growth function for a hypothesis class \mathcal{F} as follows.

$$\Pi_{\mathcal{F}}(\mathcal{F}, n) = \sup \{ |\mathcal{F}_{|x_1, \dots, x_n}| : x_1, \dots, x_n \in \mathcal{X} \}$$

Example : thresholds

What does the growth function of the class of threshold function look like ?

Well sort any given n points in ascending order, using thresholds, we can get at most $n + 1$ possible labeling on the n points. Hence $\Pi_{\mathcal{F}}(n) = n + 1$. From this we conclude that for the learning thresholds problem,

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq \sqrt{\frac{\log(n)}{n}}$$

5 Growth Function and VC dimension

Growth function is defined as,

$$\Pi(\mathcal{F}, n) = \max_{x_1, \dots, x_n} |\mathcal{F}_{|x_1, \dots, x_n}|$$

Clearly we have from the previous results on bounding minimax rates for statistical learning in terms of cardinality of growth function that :

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq \sqrt{\frac{2 \log \Pi(\mathcal{F}, n)}{n}}$$

Note that $\Pi(\mathcal{F}, n)$ is at most 2^n but it could be much smaller. In general how do we get a handle on growth function for a hypothesis class \mathcal{F} ? Is there a generic characterization of growth function of a hypothesis class ?

Definition 1. *VC dimension of a binary function class \mathcal{F} is the largest number of points $d = \text{VC}(\mathcal{F})$, such that*

$$\Pi_{\mathcal{F}}(d) = 2^d$$

If no such d exists then $\text{VC}(\mathcal{F}) = \infty$

If for any set $\{x_1, \dots, x_n\}$ we have that $|\mathcal{F}_{|x_1, \dots, x_n}| = 2^n$ then we say that such a set is shattered. Alternatively VC dimension is the size of the largest set that can be shattered by \mathcal{F} . We also define VC dimension of a class \mathcal{F} restricted to instances x_1, \dots, x_n as

$$\text{VC}(\mathcal{F}; x_1, \dots, x_n) = \max \left\{ t : \exists i_1, \dots, i_t \in [n] \text{ s.t. } \left| \mathcal{F}_{|x_{i_1}, \dots, x_{i_t}} \right| = 2^t \right\}$$

That is the size of the largest shattered subset of n . Note that for any $n \geq \text{VC}(\mathcal{F})$, $\sup_{x_1, \dots, x_n} \text{VC}(\mathcal{F}|_{x_1, \dots, x_n}) = \text{VC}(\mathcal{F})$.

1. To show $\text{VC}(\mathcal{F}) \geq d$ show that you can at least pick d points x_1, \dots, x_d that can be shattered.
2. To show that $\text{VC}(\mathcal{F}) \leq d$ show that no configuration of $d + 1$ points can be shattered.

Eg. Thresholds One point can be shattered, but two points cannot be shattered. Hence VC dimension is 1. (If we allow both threshold to right and left, VC dimension is 2).

Eg. Spheres Centered at Origin in d dimensions one point can be shattered. But even two can't be shattered. VC dimension is 1!

Eg. Half-spaces Consider the hypothesis class where all points to the left (or right) of a hyperplane in \mathbb{R}^d are marked positive and the rest negative. In 1 dimension this is threshold both to left and right. VC dimension is 2. In d dimensions, think of why $d + 1$ points can be shattered. $d + 2$ points can't be shattered. Hence VC dimension is $d + 1$.

Claim 1. VC dimension of half-spaces in \mathbb{R}^d is $d + 1$

Proof. We consider half-spaces that map vector in \mathbb{R}^d to $\{\pm 1\}$. That is

$$\mathcal{F} = \{\mathbf{x} \mapsto \text{sign}(\mathbf{f}^\top \mathbf{x} + f_0) : \mathbf{f} \in \mathbb{R}^d, f_0 \in \mathbb{R}\}$$

We prove the statement as follows.

1. $\text{VC}(\mathcal{F}) \geq d + 1$:

We can shatter the points $\mathbf{e}_1, \dots, \mathbf{e}_d, \mathbf{0}$. To see this, note that given any $y_1, \dots, y_{d+1} \in \{\pm 1\}^{d+1}$, if we consider $f \in \mathcal{F}$ given by $f_0 = y_{d+1}$ and for all $i \in [d]$, $\mathbf{f}[i] = y_i - y_{d+1}$. Hence note that, $f(\mathbf{0}) = \text{sign}(\mathbf{f}^\top \mathbf{0} + f_0) = \text{sign}(y_{d+1}) = y_{d+1}$. Also, for any $i \in [d]$, $f(\mathbf{e}_i) = \text{sign}(\mathbf{f}^\top \mathbf{e}_i + f_0) = \text{sign}(y_i - y_{d+1} + y_{d+1}) = y_i$.

2. $\text{VC}(\mathcal{F}) < d + 2$:

By Radon theorem, any set of $d + 2$ points in \mathbb{R}^d can be partitioned into two disjoint subsets whose convex hulls have a non-empty intersection. Label one of these partitions +1 and other -1. No half-space can successfully label points in the intersection.

□

Eg. Finite Hypothesis Class

Claim 2. For any binary hypothesis class \mathcal{F} ,

$$\text{VC}(\mathcal{F}) \leq \log_2 |\mathcal{F}|$$

Proof. Note that for any d , $\Pi(\mathcal{F}, d) \leq |\mathcal{F}|$. From the definition of VC dimension, we have, $\text{VC}(\mathcal{F}) = \max\{d : \Pi(\mathcal{F}, d) = 2^d\}$. Hence $2^{\text{VC}(\mathcal{F})} \leq |\mathcal{F}|$ □

Claim 3. *Learnability with binary hypothesis class \mathcal{F} implies $\text{VC}(\mathcal{F}) < \infty$.*

Proof. First note that learnability in the statistical learning framework implies learnability in the realizable PAC setting. Hence to prove the claim, it suffices to show that if a hypothesis class has infinite VC dimension, then it is not even learnable in the realizable PAC setting. To this end, assume that a hypothesis class \mathcal{F} has infinite VC dimension. This means that for any n , we can find $2n$ points x_1, \dots, x_{2n} that are shattered by \mathcal{F} . Also drawn $y_1, \dots, y_{2n} \in \{\pm 1\}$ Rademacher random variables. Let D be the uniform distribution over the $2n$ instance pairs $(x_1, y_1), \dots, (x_{2n}, y_{2n})$. Notice that since x_1, \dots, x_{2n} are shattered by \mathcal{F} , we are indeed in the realizable PAC setting for any choice of y 's. Now assume we get n input instances drawn iid from this distribution. Clearly in this sample of size n , we can at most witness n unique instances. Let us denote $J \subset [2n]$ as the indices of the $2n$ instances witnessed in the draw of n samples S . Clearly $|J| \leq n$. Hence we have,

$$\begin{aligned}
\mathcal{V}_n^{\text{PAC}}(\mathcal{F}) &\geq \sup_{x_1, \dots, x_{2n}} \inf_{\hat{y}} \mathbb{E}_{y_1, \dots, y_{2n}} \mathbb{E}_S \left[\frac{1}{2n} \sum_{j=1}^{2n} \mathbf{1}_{\{\hat{y}(x_j) \neq y_j\}} \right] \\
&= \frac{1}{2n} \sup_{x_1, \dots, x_{2n}} \inf_{\hat{y}} \mathbb{E}_{y_1, \dots, y_{2n}} \mathbb{E}_J \left[\sum_{i \in J} \mathbf{1}_{\{\hat{y}(x_i) \neq y_i\}} + \sum_{i \in [2n] \setminus J} \mathbf{1}_{\{\hat{y}(x_i) \neq y_i\}} \right] \\
&\geq \frac{1}{2n} \sup_{x_1, \dots, x_{2n}} \inf_{\hat{y}} \min_{J \subset [2n]: |J| \leq n} \mathbb{E}_{y_1, \dots, y_{2n}} \left[\sum_{i \in [2n] \setminus J} \mathbf{1}_{\{\hat{y}(x_i) \neq y_i\}} \right] \\
&= \frac{1}{4n} \min_{J \subset [2n]: |J| \leq n} |[2n] \setminus J| = \frac{n}{4n} = \frac{1}{4}
\end{aligned}$$

□