

# LECTURE 6: VC DIMENSION, FUNDAMENTAL THEOREM

Instructor: Aditya Bhaskara      Scribe: Nikhil Kala

CS 5966/6966: Theory of Machine Learning

January 27<sup>th</sup>, 2022

## Abstract

The lecture covers learnability for finite hypothesis classes. Also introduces learnability for infinite classes using Growth Functions.

## 1 CHERNOFF BOUNDS

Suppose we have a probability distribution with mean  $\mu$  and support  $[a,b]$  which is the range of the values that the distribution can take. Then if we take  $n$  iid samples  $(X_1, X_2, \dots, X_n)$ , the sample average is very close to the mean. This theorem tells us how fast the sample mean converges to the mean of the distribution.

We have:

$$\Pr\left[\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| > \epsilon\right] \leq 2 \exp^{-\frac{\epsilon^2 n}{(a-b)^2}}$$

So we can see as we increase the samples the probability drops exponentially, i.e., we converge quickly. Hence, this is also called "large deviation bound".

## 2 LEARNABILITY FOR FINITE CLASSES

If  $\mathcal{H}$  is finite hypothesis class, i.e.,  $\mathcal{H} = \{h_1, h_2, \dots, h_M\}$

**Claim:**

For any  $X$  and distribution  $D$  over it, a sample of size  $O\left(\frac{1}{\epsilon^2} \log \frac{|\mathcal{H}|}{\delta}\right)$  is representative with prob. at least  $1 - \delta$

**Proof:**

First we will look at a single  $h \in \mathcal{H}$

Now the probability that absolute difference in the sample error and risk of  $h$  is greater than  $\epsilon$  can be viewed as an application of Chernoff bound, i.e.,

$|\text{sample error}(h) - \text{risk}(h)| > \epsilon$

So to get the iid samples,

Let,

$$X_i = \begin{cases} 1 & \text{if } h \text{ is incorrect on sample } i \\ 0 & \text{if } h \text{ is correct on sample } i \end{cases}$$

So now we can say using Chernoff bound, the probability that

$| \text{sample error}(h) - \text{risk}(h) | > \epsilon$  is bounded by:

$$2e^{-\epsilon^2 m} \leq \frac{\delta}{|\mathcal{H}|}$$

Now we can use Union bound to prove that for  $m$  samples:

$\Pr[| \text{sample error}(h) - \text{risk}(h) | > \epsilon] < \delta$

To prove this we try to define bad samples, "S" for  $h$ , if

$| \text{sample error on S}(h) - \text{risk}(h) | > \epsilon$

then it is a bad sample.

Using the above bound we know that the probability that a random sample  $S$  is bad for some given  $h$  is  $\leq \frac{\delta}{|\mathcal{H}|}$  (this is true for all  $h \in \mathcal{H}$ )

This implies that the Probability that  $S$  is good for all  $h \in \mathcal{H}$  is  $\geq 1 - \delta$ .

As the union of bad probabilities is  $\frac{\delta}{|\mathcal{H}|} * m$ , as  $|\mathcal{H}| = m$ , this comes out to be  $\delta$

Thus the claim is proved.

**Note:** This essentially means that for finite hypothesis class,  $O(\log(|\mathcal{H}|))$  samples will be representative.

### 3 LEARNABILITY FOR INFINITE CLASSES

Most functions on continuous domain are infinite so we need to take into consideration infinite classes as well. Since the classes are infinite then we would need infinite samples to learn, as using the above claim will be not defined.

To overcome that we can divide the hypothesis into finitely many classes.

The logic behind this is that, if the hypothesis in a class provide the same result then we can just take the class as just one hypothesis.

Example: Threshold functions on a line

Since the sign patterns are fixed for  $m$  points as they will be  $(+++.., -++ +.., ..., ---)$   $m+1$  types of patterns we can reduce the infinite set of threshold functions to  $m+1$  classes as the hypothesis in the classes will classify the points in the same pattern.

## 4 GROWTH FUNCTIONS

From the learnability of infinite classes we reach the conclusion that, What matters isn't the number of distinct hypothesis, The number of ways in which the hypothesis classify the points of the domain matters. So to leverage that we have the growth function.

Growth function of a hypothesis class  $\mathcal{H}$  over Domain  $X$  for input size  $m$ ,

$$\tau_{\mathcal{H}}(m) = \max_{\substack{|S|=m \\ S \in X}} \{ \text{Number of distinct ways in which hypothesis in } \mathcal{H} \text{ classify } S. \}$$

**Note:** Two classification are distinct if they differ even on one example,  $+++$  is different from  $++-$

We can observe that for binary classification the growth function is  $\leq 2^m$  as those are all the possible permutations.

## 5 GROWTH FUNCTION EXAMPLES

### 1. Growth function for Linear Threshold Functions(LTFs):

$$\tau_H(m) = \max_{\substack{|S|=m \\ S \in X}} \{ \text{Number of distinct ways in which hypothesis in } \mathcal{H} \text{ classify } S. \}$$

The distinct ways are:

$$\begin{aligned} 1- &> + + + + \dots + \\ 2- &> - + + + \dots + \\ 3- &> - - + + \dots + \\ 4- &> - - - + \dots + \\ &\cdot \\ &\cdot \\ &\cdot \\ m+1- &> - - - - \dots - \end{aligned}$$

$$\tau_H(m) = m + 1$$

### 2. Growth function for intervals on number line:

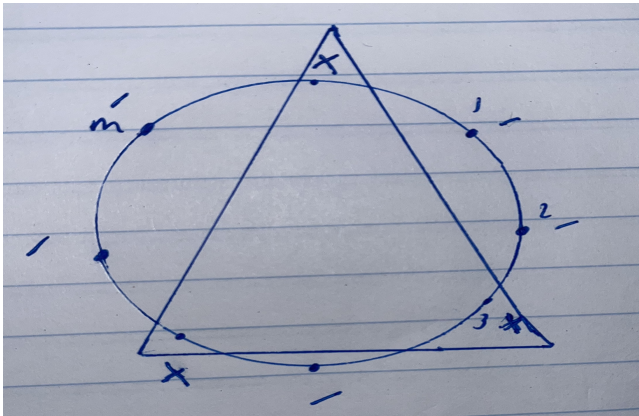
We can take interval as between any two spaces, so out of  $m+1$  spaces we select 2, so  $\binom{m+1}{2}$  which is  $O(m^2)$

$$\tau_H(m) = O(m^2)$$

### 3. Growth function for hypothesis class of convex polygons in real space:

The points inside the polygon are positive and outside are negative then to get all possible sign patterns, we can align the points equidistant on circumference of circle. Now we can take any points to construct a polygon, so we can get  $2^m$  possible combinations. (As seen in fig. for 3 points inside the polygon)

$$\tau_H(m) = O(2^m)$$



## 6 LEARNABILITY IN TERMS OF GROWTH FUNCTION

**Theorem:** Suppose  $\tau_H(m)$  is an upper bound on the total number of distinct “classifications” (or “sign patterns”) possible for any sample of size  $m$ . Then for any  $X, D$ , if we take a sample  $S$  of size  $m$ , we have, with prob.  $1-\delta$

$$\underbrace{\sup_{h \in H} |err(h, S) - err(h, D)|}_{\epsilon} \leq \frac{4 + \sqrt{\log \tau_H(2m)}}{\delta \sqrt{2m}}$$

So to make this  $\epsilon$ -representative, so

$$\frac{4 + \sqrt{\log \tau_H(2m)}}{\delta \sqrt{2m}} = \epsilon$$

Evaluating this for  $\tau_H(m) = m + 1$

We get,

$$m = \frac{4 \log \frac{1}{\epsilon \delta}}{\epsilon^2 \delta^2}$$