

# Computational and Statistical Learning Theory

## Problem set 1

Due: April 15th

Please send your solutions to `learning-submissions@ttic.edu`

### Notation:

- Input space:  $\mathcal{X}$
- Label space:  $\mathcal{Y} = \{\pm 1\}$
- Sample:  $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X}$
- Hypothesis Class:  $\mathcal{H}$
- Risk:  $L_{\mathcal{D}}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbf{1}_{h(x) \neq y}]$
- Empirical Risk:  $L_S(h) = \frac{1}{m} \sum_{(x,y) \in S} \mathbf{1}_{h(x) \neq y}$

### 1. Binomial Tail Bounds:

$\{0, 1\}$ -valued random variables  $X_1, \dots, X_n$  are drawn independently each from Bernoulli distribution with parameter  $p = 0.1$ . Define  $P_n := \mathbb{P}(\frac{1}{n} \sum_{i=1}^n X_i \leq 0.2)$ .

- For  $n = 1$  to 30 calculate and plot the below in the same plot (see [1, section 6.1] for definition of Hoeffding and Bernstein inequalities):
  - Exact value of  $P_n$  (binomial distribution).
  - Normal approximation for  $P_n$ .
  - Hoeffding inequality bound on  $P_n$ .
  - Bernstein inequality bound on  $P_n$ .
- For  $n = 30$  to 300 calculate and plot the below in the same plot :
  - Normal approximation for  $P_n$ .
  - Hoeffding inequality bound on  $P_n$ .
  - Bernstein inequality bound on  $P_n$ .

### 2. VC Bound:

Given a set  $C = \{x_1, \dots, x_m\}$  let  $\mathcal{H}_{x_1, \dots, x_m} = \{(h(x_1), \dots, h(x_m)) \in \{\pm 1\}^m : h \in \mathcal{H}\}$ . Recall that we say that such a set is *shattered* by  $\mathcal{H}$  if  $|\mathcal{H}_{x_1, \dots, x_m}| = 2^m$ , and that the VC

dimension of  $\mathcal{H}$  is the size of the largest sample set that can be shattered. Also recall that the *growth function* of the hypothesis class  $\mathcal{H}$  is given by:

$$\Gamma_{\mathcal{H}}(m) = \sup_{x_1, \dots, x_m} |\mathcal{H}_{x_1, \dots, x_m}|. \quad (1)$$

That is, we can also define the VC dimension as the largest  $m$  for which  $\Gamma_{\mathcal{H}}(m) = 2^m$ .

In this problem, we will see how to obtain a uniform convergence guarantee, bounding the differences between empirical and expected errors, in terms of the growth function. We already know how uniform convergence ensures learning guarantees for ERM, and so we can obtain learning guarantees in terms of the growth function, and thus using 1 also in terms of the VC dimension.

The growth function bounds the number of behaviors of the hypothesis class on a finite number of points, while the expected error depends on the behavior on all the points. To overcome this, we will first bound the difference between two the empirical error on two different samples, and then show that the difference between the expected and empirical errors cannot be much larger than the difference between two different empirical errors. This technique is called “systematization”: we are introducing a “ghost sample” of another  $m$  points sampled i.i.d. from the source distribution to stand in for the expected error and make the problem symmetric.

In problem 5, you will prove that if  $\mathcal{H}$  has VC dimension  $d$ , then for any  $m$ ,

$$\Gamma_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d \quad (2)$$

We will now see how to use systematization to obtain a learning guarantee that depends on the growth function, and hence on the VC dimension.

- (a) For any sequence of  $2m$  points  $S = (z_1, \dots, z_m, z'_1, \dots, z'_m)$ , consider  $m$  i.i.d. uniform random signs  $s_1, \dots, s_m$  which define the samples  $S_1, S_2$  in the following way: for each  $i = 1 \dots m$ , if  $s_i = 1$  then  $z_i \in S_1$  and  $z'_i \in S_2$ , otherwise (if  $s_i = -1$ ) then  $z_i \in S_2$  and  $z'_i \in S_1$ ; i.e. the variables  $s_1, \dots, s_m$  specify how to “deal” the  $2m$  points into the two sets  $S_1$  and  $S_2$ . Now, for any sequence  $S$  of  $2m$  points, and any hypothesis  $h$ , with  $\ell(h, z) \in \{0, 1\}$ , prove that with probability  $\geq 1 - \delta$  over the separation to  $S_1, S_2$ :

$$|L_{S_1}(h) - L_{S_2}(h)| \leq \sqrt{f(\delta)/m} \quad (3)$$

Hint: Write  $L_{S_1}(h) - L_{S_2}(h) = \frac{1}{m} \sum_{i=1}^m s_i (\ell(h, z_i) - \ell(h, z'_i))$

- (b) For any sequence  $S$  of  $2m$  points as above, prove that with probability  $\geq 1 - \delta$  over the separation to  $S_1, S_2$ , for every  $h \in \mathcal{H}$ :

$$|L_{S_1}(h) - L_{S_2}(h)| \leq \sqrt{f(\delta, \Gamma_{\mathcal{H}}(2m))/m} \quad (4)$$

and conclude that for any hypothesis class  $\mathcal{H}$ , any source distribution  $\mathcal{D}$  and any sample size  $m$ , with probability at least  $1 - \delta$  over  $S_1, S_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}^m$  (i.e. each sample is drawn independently from  $\mathcal{D}^m$ ), for all  $h \in \mathcal{H}$ ,

$$|L_{S_1}(h) - L_{S_2}(h)| \leq \sqrt{f(\delta, \Gamma_{\mathcal{H}}(2m))/m} \quad (5)$$

- (c) (Optional) Prove the symmetrization lemma; i.e. prove that for any hypothesis class  $\mathcal{H}$ , any source distribution  $\mathcal{D}$ , any  $\epsilon > 0$  and any sample size  $m$  such that  $m \geq \frac{1}{2\epsilon^2}$ :

$$\mathbb{P}_{S \sim \mathcal{D}^m} (\exists_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| > 2\epsilon) \leq 2\mathbb{P}_{S, S' \sim \mathcal{D}^m} (\exists_{h \in \mathcal{H}} |L_S(h) - L_{S'}(h)| > \epsilon) \quad (6)$$

Hint: Prove the above inequality is two steps:

- i. Show that for any  $S$  and  $S'$ , the following inequality holds:

$$\mathbf{1}_{\exists_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| > 2\epsilon} \mathbf{1}_{\exists_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_{S'}(h)| < \epsilon} \leq \mathbf{1}_{\exists_{h \in \mathcal{H}} |L_S(h) - L_{S'}(h)| > \epsilon}$$

- ii. Take the expectation with respect to  $S'$  and use Chebychev's inequality to bound  $\mathbb{P}(\exists_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_{S'}(h)| < \epsilon)$  and then take the expectation with respect to  $S$ . Chebychev's inequality bounds the probability of deviation from the expected value:

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \frac{\text{Var}(X)}{t^2}$$

- (d) Use the symmetrization lemma and part (c) above, to prove that, with probability  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^m$ , for all  $h \in \mathcal{H}$ :

$$|L_S(h) - L_{\mathcal{D}}(h)| < \sqrt{f(\delta, \Gamma_{\mathcal{H}}(2m))/m} \quad (7)$$

and conclude that if  $\text{VC-dim}(\mathcal{H}) \leq D$  then:

$$L_{\mathcal{D}}(\hat{h}) \leq L_{\mathcal{D}}(h^*) + \sqrt{f(\delta, D \log(2em/D))/m} \quad (8)$$

- (e) Conclude that  $m = O(D \log(1/\epsilon)/\epsilon^2)$  samples are enough to ensure that with probability  $\geq 1 - \delta$ ,  $L_{\mathcal{D}}(\hat{h}) < L_{\mathcal{D}}(h^*) + \epsilon$ . Write down an explicit bound (without big-O notation, though the constants need not be the tightest possible).

Hint: start with the expression for  $m$ , plug it into the r.h.s. of part (d) above, and verify that the r.h.s is less than  $\epsilon$ .

### 3. VC Dimension of Sparse Linear Classifiers Using Concept Class Union Arguments:

Fix an instance space  $\mathcal{X}$ . Let  $\{\mathcal{H}_i\}$  be a family of concept classes over  $\mathcal{X}$ , where  $i$  ranges from 1 to some integer  $r$ . Define the concept  $\mathcal{H}$  to be the union  $\cup_{i=1 \dots r} \mathcal{H}_i$ .

- (a) Give an upper bound for the growth function  $\Gamma_m(\mathcal{H})$  in terms of the growth functions  $\Gamma_m(\mathcal{H}_i)$  for  $i = 1 \dots r$ .
- (b) If the VC dimension of all  $\mathcal{H}_i$  are bounded by  $D$ , i.e.  $\text{VCdim}(\mathcal{H}_i) \leq D$ , give a bound on the growth function  $\Gamma_m(\mathcal{H})$  in terms of  $m$ ,  $r$  and  $D$ .

Hint: Use Sauer's lemma, which you will prove in problem 5.

- (c) From (b) conclude that  $\text{VCdim}(\mathcal{H}) = O(\max(D, \log(r) + D \log(\log(r)/D)))$ . Give an exact bound, without the  $O(\cdot)$  notation.

For the remainder of the question, consider a feature mapping  $\phi : \mathcal{X} \mapsto \mathbb{R}^d$ . For a vector  $w \in \mathbb{R}^d$ , we say that  $w$  is  $k$ -sparse if the number of coordinates  $j$  for which  $w(j) \neq 0$  is at most  $k$ . We Define the following concept class:

$$\mathcal{H}^{(k)} := \{h_w(x) = \text{sign}(\langle w, \phi(x) \rangle) \mid \text{s.t. } w \text{ is } k\text{-sparse}\} .$$

- (d) Using (c) above, show that the VC dimension of  $\mathcal{H}^{(k)}$  is at most  $O(k \log(d/k))$ .
- (e) Show how to shatter a subset of size  $\Omega(\log d)$  with respect to  $\mathcal{H}^{(1)}$ , establishing tight upper and lower bounds on the VC dimension of  $\mathcal{H}^{(1)}$ .
- (f) (Optional) Show how to shatter a subset of size  $\Omega(k \log(d/k))$  with respect to  $\mathcal{H}^{(k)}$ , establishing tight upper and lower bounds on the VC dimension of  $\mathcal{H}^{(k)}$ .

#### 4. The Order of Quantifiers in Statistical No Free Lunch:

Notice the order of quantifiers in the converse of the fundamental theorem: for any hypothesis class and any learning rule, there exists a distribution on which we need at least  $\text{VCdim}(\mathcal{H})/2$  samples to get small error with high probability. That is, some distributions might be easy while others might be hard, and the “hard distribution” depends on the learning rule. Can these quantifiers be reversed?

- (a) Show that for any hypothesis class and any distribution, there exists a learning rule that with a very small number of samples returns a predictor that is as good as the best predictor in the class, i.e. obtains generalization error  $\inf_{h \in H} L_{\mathcal{D}}(h)$ .
- (b) Consider  $\mathcal{X} = \mathbb{R}$ , and the hypothesis class:

$$\mathcal{H} = \{[x \in C] \mid C \text{ is a finite subset of } \mathbb{R}_-\} \cup \{[a < x < b] \mid 0 < a < b < \infty\}$$

- i. What is the VC dimension of  $\mathcal{H}$ ?
- ii. For distributions supported only on  $\mathbb{R}_+$ , how many samples are required in order for ERM to obtain error  $\inf_{h \in H} L_{\mathcal{D}}(h) + \epsilon$  with probability at least 0.99 ?

#### 5. (Optional) Sauer’s Lemma:

The aim of this exercise is to prove the “Sauer’s Lemma”: if  $\mathcal{H}$  has VC dimension  $d$ , then for any  $m$ ,

$$\Gamma_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}. \quad (9)$$

In order to prove (9), we will actually prove the following statement: for any set  $C = \{x_1, \dots, x_m\}$ :

$$|\mathcal{H}_C| \leq |\{B \subset C : B \text{ is shattered by } \mathcal{H}\}| \quad (10)$$

That is, the number of possible labeling of a  $C$  is bounded by the number of different subsets of  $C$  that can be shattered.

We (i.e. you) will prove (10) by induction.

- (a) Establish that (10) holds for  $C = \emptyset$  (the empty set).
- (b) For any set  $C$  and any point  $x' \notin C$ , assume (10) holds for  $C$  and for any hypothesis class, and prove that (10) holds for  $C' = C \cup \{x'\}$  and any hypothesis class. To this end, for any hypothesis class  $\mathcal{H}$ , write  $\mathcal{H} = \mathcal{H}^- \cup \mathcal{H}^+$  where:

$$\mathcal{H}^+ = \{h \in \mathcal{H} : h(x') = +1\}$$

$$\mathcal{H}^- = \{h \in \mathcal{H} : h(x') = -1\}$$

- i. Prove that  $|\mathcal{H}_{C'}| = |\mathcal{H}_C^+| + |\mathcal{H}_C^-|$ .
- ii. Prove that (10) holds for  $C$  and  $\mathcal{H}$  by applying (10) to each of the two terms on the right-hand-side above.

We can now conclude that (10) holds for any (finite)  $C$  and any  $\mathcal{H}$ .

- (c) Use (10) to establish (9).
- (d) For  $d \leq m$ , prove that  $\sum_{i=0}^d \binom{m}{i} \leq m^d + 1$ .
- (e) Prove the tighter bound:  $\sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d$ .

## References

- [1] O. Bousquet, S. Boucheron, and G. Lugosi. *Introduction to statistical learning theory*. Advanced Lectures on Machine Learning, pp. 169-207. Springer Berlin Heidelberg, 2004.