

# Computational Learning Theory Lecture Notes

Martin Citoler-Saumell

CIS625 Spring 2017

These are notes from [Prof. Michael Kearns](#)'s CIS625. This course is currently being taught at the University of Pennsylvania and attempts to provide algorithmic, complexity-theoretic and probabilistic foundations to modern machine learning and related topics. You can find out more at the [course website](#). A word of caution, it is likely that some typos and/or errors have been introduced during the taking of these notes. If you find any, please let me know at [martinci@math.upenn.edu](mailto:martinci@math.upenn.edu).

## 1 Lecture 1: 2017.01.23

### 1.1 Course Outline/Description

#### 1.1.1 Formal Models of ML

- Assumptions about data generation process.
- Assumptions about what algorithms "knows".
- Sources of info that the algorithms has
- Criteria/objective of learning is.
- Restrictions on algorithms.

#### 1.1.2 Examples of Models

- "PAC" model (in first 1/2 of the term).
- Statistical learning theory.
- "no-regrets" learning models.
- Reinforcement learning.
- ML & Differential privacy.
- "Fairness" in ML

## 1.2 A Rectangle Learning Problem

Suppose you are trying to teach an alien friend the "shape" of humans in terms of abstract descriptions like "medium build", "athletic", etc. We are going to assume that each one of these descriptions represents a rectangular region on the height-weight plane but we are not aware of the exact dimensions. The only thing we are able to tell the alien is whether a particular individual is medium built or not. i.e. we can only label examples.

- target rectangle  $R$ , the *true* notion of "medium built".
- hypothesis rectangle  $\hat{R}$ .

*Remark 1.1.* Note that the assumption that the classifier function is a rectangle is rather strong. There is always this trade-off to be able to actually compute things. From a Bayesian point of view, "we always need a prior".

Given a data cloud of examples, a reasonable hypothesis rectangle could be the tightest fit rectangle. However, this choice ignores the negative examples so it seems that that we are throwing away information. In a sense, this rectangle would be the least likely.

- Assume  $\langle x_1, x_2 \rangle$  pairs of height-weight are i.i.d. from an arbitrary unknown probability distribution,  $D$ .

We want to be able to evaluate how our hypothesis rectangle is performing. We want bounds on the classification error, which can be thought as the size of the symmetric difference between  $R$  and  $\hat{R}$

$$D[R \triangle \hat{R}] \equiv \mathbb{P}_D[R(\vec{x}) \neq \hat{R}(\vec{x})]. \quad (1.1)$$

**Theorem 1.2.** *Given  $\varepsilon, \delta > 0$ , there is some integer  $N$  such that if we have more than  $N$  training examples,<sup>1</sup> we have*

$$D[R \triangle \hat{R}] \equiv \mathbb{P}_D[R(\vec{x}) \neq \hat{R}(\vec{x})] < \varepsilon, \quad (1.2)$$

*with probability at least  $1 - \delta$ .*

*Proof.* First of all, note that using tightest fit, the hypothesis rectangle is always contained inside the target rectangle. Now, for each side of the target rectangle we may draw inward strips in such a way that each strip has  $\mathbb{P}_D[\text{Strip}] < \frac{\varepsilon}{4}$ . If the training set has a positive example in each of these four strips, then the inequality above is satisfied because the boundary of the hypothesis rectangle would be contained in the union of the strips. Next we need to deal with the required sample size to obtain this result with some certainty. Let  $m$  denote the sample size, since the distribution is i.i.d., we have

$$\begin{aligned} \mathbb{P}_D[\text{miss a specific strip } m \text{ times}] &= \left(1 - \frac{\varepsilon}{4}\right)^m, \\ \mathbb{P}_D[\text{miss any of the strips } m \text{ times}] &\geq 4 \left(1 - \frac{\varepsilon}{4}\right)^m. \end{aligned}$$

---

<sup>1</sup>This the same as saying that sample size is at least  $N$ .

By the discussion above, the last inequality implies

$$\mathbb{P}_{D^m}[D[R\Delta\hat{R}] \geq \varepsilon] \leq 4 \left(1 - \frac{\varepsilon}{4}\right)^m,$$

which can be chosen arbitrarily small for big enough  $m$ . One can obtain  $N \geq \frac{4}{\varepsilon} \ln\left(\frac{4}{\delta}\right)$ .  $\square$

*Remark 1.3.* This proof generalizes to  $d$ -dimensional rectangles. We only need to replace 4 with  $2d$ , the number of  $(d-1)$ -faces. We can also try to incorporate noisy data, where the labels have some probability of being wrong.

### 1.3 A More General Model

- Input/instance/feature space,  $X$ . (e.g.  $\mathbb{R}^2$  in the example above)
- Concept/classifier/boolean function,  $C : X \rightarrow \{0, 1\}$  or we can also think about it as an indicator function of the positive examples or the subset of positive examples.
- Concept class/target class,  $\mathcal{C} \subset \mathcal{P}(X)$ , the admissible concepts/classifiers. (e.g. all rectangles in  $\mathbb{R}^2$  in the example above)
- target concept,  $c \in \mathcal{C}$ . (e.g. target rectangle in the example above)
- Input distribution,  $D$  over  $X$  (arbitrary & unknown)
- Learning algorithm given access to examples of the form:  $\langle \vec{x}, y \rangle$  where  $\vec{x}$  is i.i.d. drawn from  $D$  and  $c(\vec{x}) = y$ .

**Definition 1.4** (PAC Learning). We say that a class of functions over  $X$ ,  $\mathcal{C}$ , is *Probably Approximately Correct (PAC) learnable* if there exists an algorithm,  $L$ , such that given any  $c$  in  $\mathcal{C}$ ,  $D$  a distribution over  $X$  and  $\varepsilon, \delta > 0$ ,  $L$  with these parameters and random inputs  $\vec{x}$ 's satisfies:

- (Learning) With probability  $\geq 1 - \delta$ ,  $L$  outputs a hypothesis,  $h$  in  $\mathcal{C}$  such that  $D[h\Delta c] < \varepsilon$ , i. e.

$$Err(h) := \mathbb{P}_{x \sim D}[h(x) \neq c(x)] \leq \varepsilon. \quad (1.3)$$

- (Efficient)  $L$  runs in time/sample  $poly\left(\frac{1}{\varepsilon}, \frac{1}{\delta}, \text{dimension}\right)$ .

## 2 Lecture 2: 2017.01.30

*Remark 2.1* (PAC Learning). Fixing the class  $\mathcal{C}$  is a strong assumption, it is the prior you are assuming about the true behavior of the data. For example, when you fit a linear regression, you are assuming that there is a true linear relation in the data.

In contrast, the assumption on the distribution over  $X$  is fairly general.

**Theorem 2.2.** *The class of rectangles over  $\mathbb{R}^2$  from example above is PAC learnable (with sample size  $m \sim \frac{1}{\varepsilon} \ln\left(\frac{1}{\delta}\right)$ ).*

## 2.1 PAC learning boolean conjunctions

In the following we are going to see an example of problem that is PAC learnable.

- $X = \{0, 1\}^n$
- Class  $\mathcal{C}$  = all conjunctions over  $x_1, \dots, x_n$ .  $|\mathcal{C}| = 3^n$   
E.g.: If  $c = x_1 \wedge \neg x_3 \wedge \neg x_{11} \wedge x_{26} \dots$ ,

$$c(\vec{x}) = 1 \iff x_1 = 1, x_3 = 0, x_{11} = 0, x_{26} = 1, \dots \quad (2.1)$$

- $D$  over  $\{0, 1\}^n$ .

### 2.1.1 Algorithm for monotone case (i.e. no $\neg$ 's)

In this case we are trying to fit something like  $c = x_1 \wedge x_5 \wedge x_{13} \dots$  i.e. given some examples of sequences of bits and the result of  $c$  on them, we are trying to guess what the conjunction  $c$  actually is. We can use the following algorithm:

- $h \leftarrow x_1 \wedge x_2 \wedge x_3 \wedge \dots \wedge x_n$ , start with the conjunction of all the variables.
- For each positive example,  $\langle \vec{x}, 1 \rangle$ , delete any variable in  $h$  such that  $x_i = 0$ .  
This method ensures that the positives of  $h$  is a subset of the true  $c$ .

### Analysis

Let  $p_i$  denote the probability we delete  $x_i$  from  $h$  in a single draw. In other words,

$$p_i = \mathbb{P}_{\vec{X} \sim D}[c(\vec{x}) = 1, x_i = 0]. \quad (2.2)$$

Then we have an a priori bound on error:  $Err(h) \leq \sum_{x_i \in h} p_i$ . We can make a distinction between *bad* and *good* indices. An index,  $i$ , is bad if  $p_i \geq \frac{\varepsilon}{n}$  and good otherwise. Note that if  $h$  contains no bad indices, then we have

$$Err(h) \leq \sum_{x_i \in h} p_i \leq n \left( \frac{\varepsilon}{n} \right). \quad (2.3)$$

Let's fix some index,  $i$ , such that  $p_i \geq \frac{\varepsilon}{n}$ . i.e. a bad index. We have that

$$\mathbb{P}[x_i \text{ "survives" } m \text{ random samples}] = (1 - p_i)^m \quad (iid) \quad (2.4)$$

$$\leq \left( 1 - \frac{\varepsilon}{n} \right)^m \quad (2.5)$$

$$\mathbb{P}[\text{any bad } i \text{ "survives" } m \text{ random samples}] \leq n \left( 1 - \frac{\varepsilon}{n} \right)^m. \quad (2.6)$$

If we want the right-hand-side of the last inequality to be less than some  $\delta > 0$ , we end up with  $m \geq \frac{n}{\varepsilon} \ln \left( \frac{n}{\delta} \right)$ . In other words, we just proved the following theorem

**Theorem 2.3.** *Conjunctions over  $\{0, 1\}^n$  are PAC learnable with sample size  $m \geq \frac{n}{\epsilon} \ln \left( \frac{n}{\delta} \right)$ .*

*Remark 2.4.* An analogous argument proves this theorem for conjunctions that are not necessarily monotone. The only difference is that we have to keep track the extra  $\neg$  variables.

*Remark 2.5.* We can identify some pattern for this kind of analysis. We identify some “bad” things that may happen and then prove that the probability of them happening decreases fast when we increase the number of samples seen.

## 2.2 Hardness of PAC learning 3-term DNF

Now we are going to see an example of non PAC learnable problem. However, we will be able to slightly modify it and achieve PAC learning. This motivates an better definition of PAC learnable.

- Input space  $X = \{0, 1\}^n$
- Class  $\mathcal{C} =$  all disjunctions of three conjunctions.  $|\mathcal{C}| = 3^{3n}$   
 E.g.: If  $c = T_1 \vee T_2 \vee T_3$  where  $T_i$  is a conjunction over  $X$ .  
 $c = (x_1 \wedge x_7 \wedge \neg x_8) \vee (x_2 \wedge x_5 \wedge \neg x_7 \wedge \neg x_8) \vee (x_6 \wedge x_{12})$ .

To see that this problem is *hard*, we prove that the graph 3-coloring problem reduces to the 3-term DNF problem.

### Graph 3-coloring to PAC learning 3-term DNF

Suppose that you have a 3-colorable (undirected) graph  $G$ . That is, a graph such that we can color the vertices with 3 colors in such a way that there are no edges between vertices of the same color. We see an example of how to transform such a graph into a set of labeled examples for the PAC learning 3-term DNF.

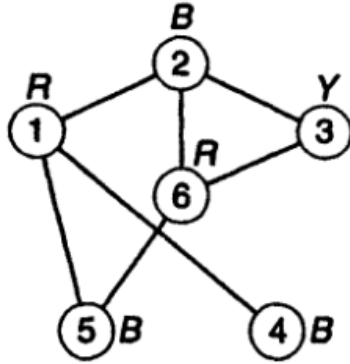


Figure 1: A graph with a 3-coloring.

First, for the  $i$ -th node we create a positive example that is represented by the vector,  $v(i)$ , with a 0 in the  $i$ -th entry and 1's everywhere else. E.g. from node one we have  $\langle 011111, + \rangle$ . Then we create a negative example from each edge that is represented by a vector,  $e(i, j)$ ,

of 1's except for 0's at the positions that determine the edge. E.g. from the edge connecting 2 and 3 we obtain  $\langle 100111, - \rangle$ . Next, the coloring can be used to define a 3-term DNF in the following way. Given a color, we define  $T_{color}$  as the conjunction of the variables/vertices that are *not* of that color. In this example we have

$$T_R = x_2 \wedge x_3 \wedge x_4 \wedge x_5, \quad (2.7)$$

$$T_B = x_1 \wedge x_3 \wedge x_6, \quad (2.8)$$

$$T_Y = x_1 \wedge x_2 \wedge x_4 \wedge x_5 \wedge x_6. \quad (2.9)$$

**Claim 2.6.**  $G$  is 3-colorable  $\iff$  there is a 3-term DNF consistent with the labeled sample above.

*Proof.* We have just seen how to obtain a 3-term DNF from a 3-colorable graph. We only need to check that it is consistent with the sample. By construction, all the positive examples satisfy  $T_{color}$  where color is the vertex's color, since the only 0 in the vector is in the position that is dropped. Similarly, it follows that the examples coming from the edges do not satisfy any of the  $T$ 's. In any edge there are two colors corresponding to its vertices. The extra 0 in the example ensures that the  $T$ 's from those colors are not satisfied. Finally, the  $T$  of the remaining color cannot be satisfied because both vertices are included in the conjunction and they are both 0 in the example.

Conversely, given a graph and the labeled examples as above, if there is a consistent 3-term DNF we can find a coloring in the following way. Label each term of the formula with a color, say  $T_R \vee T_Y \vee T_B$ , and remember the order of the labels. Then, we define the color of vertex  $i$  (corresponding to vector  $1 \dots 101 \dots 1$ ) as the label of the first formula that is satisfied by  $v(i)$ . Since the formula is consistent with the sample, every vertex must be actually colored. We only need to argue that this is a valid coloring. Suppose to the contrary that  $i$  and  $j$  are to vertices that are connected by an edge and have the same color. This means that both  $v(i)$  and  $v(j)$  satisfy  $T_{C_0}$ . However, we also have  $v(i) \& v(j) = e(i, j)$  where  $\&$  denotes the bit-wise and operation and it follows that  $e(i, j)$  satisfies  $T_{C_0}$ , a contradiction with the consistency of the formula since edges are negative examples.  $\square$

This concludes the argument that finding a coloring of a graph is the same as producing a consistent 3-term DNF. We are only left with the computational aspects. Namely, given the labeled sample associated to a graph, we need to find a way to feed it to a PAC learning algorithm. First we need a distribution over vectors of bits. For this we can just sample the examples uniformly. Finally, to ensure consistency we choose  $\varepsilon$  any quantity less than  $\frac{1}{\#examples}^2$ . This way the algorithm cannot make any mistakes is forced to be consistent. The  $\delta$  can be arbitrary. In conclusion, if the 3-term DNF problem were PAC learnable, we could solve the coloring problem in random polynomial time. Another way to say this is in the form of the following theorem.

**Theorem 2.7.** *If 3-term DNF are PAC learnable, then  $NP = RP$ .*

Of course, it is strongly believed that  $RP \subsetneq NP$  so this is rather strong evidence against the easiness of the 3-term DNF problem.

---

<sup>2</sup>This epsilon is allowed because  $\frac{1}{\varepsilon}$  is polynomial in the size of input. This is not true in general.

*Remark 2.8.* The upshot is that a slight generalization of the conjunction problem, for which we have a randomized polynomial time solution, is almost assured to not be PAC learnable ( unless  $NP = RP$ )

### 3 Lecture 3: 2017.02.06

Recall the definition of a class being PAC learnable.

**Definition 3.1** (PAC Learning). A class  $\mathcal{C}$  is *Probably Approximately Correct (PAC) learnable* if there exists an algorithm,  $L$ , such that:

$$\forall c \in \mathcal{C}, \quad \forall D \text{ over } X, \quad \forall \varepsilon, \delta > 0 \quad (3.1)$$

- (Learning) With probability  $\geq 1 - \delta$ ,  $L$  outputs a hypothesis,  $h$  in  $\mathcal{C}$  such that  $D[h \Delta c] < \varepsilon$ , i.e. we have error at most  $\varepsilon$

$$Err(h) := \mathbb{P}_{x \sim D}[h(x) \neq c(x)] \leq \varepsilon. \quad (3.2)$$

- (Efficient)  $L$  runs in time/sample  $\text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, n)$ . *Samples & computation.*

*Remark 3.2.* Usually when we talk about computation time we are in the realm of complexity theory and we talk about samples we are really asking statistics/information-theory questions about what sample size do we need to be able to draw some conclusion.

#### 3.1 PAC learning 3-term DNF by 3-CNF

In the following we see how to overcome the intractability of the 3-DNF PAC learning by using a different representation. It amounts to expanding the input space and the class we are learning.

**Definition 3.3.** A 3-CNF is a conjunction of disjunctions of length three.

Given a 3-DNF, we can rewrite it as a 3-CNF in the following way:

$$T_1 \vee T_2 \vee T_3 \equiv \bigwedge_{\substack{u \in T_1 \\ v \in T_2 \\ w \in T_3}} (u \vee v \vee w), \quad (3.3)$$

for every assignment of  $x_1, \dots, x_n$ , both sides evaluate to the same boolean value. Notice that the length of the 3-CNF can be much bigger (but still polynomial): 3-DNF is at most as  $3n$  but 3-CNF could be up to  $n^3$ . In a sense, this corresponds to feature generation.

*Remark 3.4.* Notice that the reverse is NOT true. Given a 3-CNF it might not be representable as a 3-DNF.

Another important point to notice is that after the transformation into 3-CNF we are changing the distribution of the initial input space. But the definition of PAC learnable allows for *any* distribution, so we are fine.

The upshot here is that we can learn 3-CNF by 3-CNF but this problem contains the intractable problem of learning 3-DNF by 3-DNF. The trick is that we have a bigger solution space so the 3-CNF algorithm is fed a 3-DNF, it has the option to output something outside the 3-DNF class, namely a 3-CNF.

**Theorem 3.5.** *3-CNF is PAC-learnable and 3-DNF. Further, by the discussion above, 3-DNF is learnable “by” 3-CNF.*

This motivates the more general definition for PAC-learnable that takes into account the solution class.

**Definition 3.6** (PAC Learning). A class  $\mathcal{C}$  is *Probably Approximately Correct (PAC) learnable* by  $\mathcal{C} \subset \mathcal{H}$  if there exists an algorithm,  $L$ , such that:

$$\forall c \in \mathcal{C}, \quad \forall D \text{ over } X, \quad \forall \varepsilon, \delta > 0 \quad (3.4)$$

- (Learning) With probability  $\geq 1 - \delta$ ,  $L$  outputs a hypothesis,  $\underline{h} \in \mathcal{H}$  such that  $D[h\Delta c] < \varepsilon$ , i.e. we have error at most  $\varepsilon$

$$Err(h) := \mathbb{P}_{x \sim D}[h(x) \neq c(x)] \leq \varepsilon. \quad (3.5)$$

- (Efficient)  $L$  runs in time/sample  $\text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta}, n)$ . *Samples & computation.*

*Remark 3.7.*  $\mathcal{C}$  is usually called the target class and  $\mathcal{H}$  the hypothesis class.

## 3.2 Consistency implies PAC-learnable

- Suppose we have some target and hypothesis classes,  $\mathcal{C} \subset \mathcal{H}$ .
- Let  $A$  be a *consistent* algorithm:
  - i) Given any finite sample,  $S = \{\langle x_1, y_1 \rangle, \dots, \langle x_m, y_m \rangle\}$ , where for all  $i$  we have  $y_i = c(x_i)$  for some  $c \in \mathcal{C}$ .
  - ii)  $A$  outputs  $h \in \mathcal{H}$  such that  $h(x_i) = y_i = c(x_i)$  for all  $i$ .

**Theorem 3.8.** *For any finite  $\mathcal{H}$ , a consistent algorithm for  $\mathcal{H}$  is a PAC-learnable algorithm.*

*Proof.* We call a hypothesis,  $h \in \mathcal{H}$ ,  $\varepsilon$ -bad if  $Err(h) > \varepsilon$ . Now we generate a size- $m$   $S$  according to  $Ex(c, D)$ , which is the subroutine that generates samples from  $D$ . For any fixed  $\varepsilon$ -bad  $h$ , we have an upper bound on the probability of  $h$  being consistent with  $S$

$$\mathbb{P}_S[h \text{ consistent with } S] \leq (1 - \varepsilon)^m. \quad (3.6)$$

Therefore,  $\mathbb{P}_S[\exists h \in \mathcal{H} \text{ that is both } \varepsilon\text{-bad and consistent with } S] \leq |\mathcal{H}|(1 - \varepsilon)^m$ . Now we choose  $\delta > 0$  such that  $|\mathcal{H}|(1 - \varepsilon)^m < \delta$ . We can conclude that as long as  $m \geq \frac{\text{const}}{\varepsilon} \ln \frac{|\mathcal{H}|}{\delta}$  the PAC learning definition will be satisfied.  $\square$



*Remark 3.9.*  $|\mathcal{H}|(1 - \varepsilon)^m \leq e^{\ln|\mathcal{H}| + m \ln(1 - \varepsilon)} \approx e^{\ln|\mathcal{H}| - c_0 \varepsilon m} \rightarrow 0$  as  $m \rightarrow \infty$ . A key point of this is that we can let the complexity of the hypothesis to grow as the sample size gets bigger and keeping this quantity under control. That is, the more data you have, the more complex the model you can train without over-fitting too much. If  $\mathcal{H}_m$  is the hypothesis for data of size  $m$ , we only need  $\ln|\mathcal{H}_m| \leq c \cdot m^\beta$ , for some  $\beta < 1$  and  $c = c(\dim)$ . From here we obtain  $m \geq \left(\frac{c}{\varepsilon}\right)^{\frac{1}{1-\beta}}$ .

Next we want to deal with the case of a possibly infinite hypothesis class  $\mathcal{H}$ . The first approach could be to try and force feed a discretization of the hypothesis into the finite case one but this might not work because it depends on the interaction between the distribution of the data and the chosen discretization. We want something better.

- Class  $\mathcal{H}$  over  $X$ .
- $h \in \mathcal{H}$  as functions  $h(x) \in \{0, 1\}$ ; or as sets  $h \subset X$ .
- Let  $S = \{x_1, \dots, x_m\}$  be an ordered subset of  $X$ .
- $\Pi_{\mathcal{H}}(S) = \{\langle h(x_1), \dots, h(x_m) \rangle : h \in \mathcal{H}\} \subseteq \{0, 1\}^m$ ,  $|\Pi_{\mathcal{H}}(S)| \leq 2^m$ .

*Remark 3.10.* If the inclusion above is saturated then it means that our hypothesis can classify ANY labeling of that size. That is bad for learning because it is saying that there is no structure in the data.

*Remark 3.11.* For each variable in  $S$ , say  $x_i$ , it induces a partition on the class  $\mathcal{H}$  according on what is the value of  $h(x_i)$  for  $h \in \mathcal{H}$ . Then, one can think about  $\Pi_{\mathcal{H}}(S)$  as the collection of all the the possible labelings with concepts in the class  $\mathcal{H}$ .

**Definition 3.12.** We say that  $S$  is *shattered* if the equality case holds,  $\Pi_{\mathcal{H}}(S) = \{0, 1\}^m$ , or equivalently,  $|\Pi_{\mathcal{H}}(S)| = 2^m$ .

**Definition 3.13.** The *Vapnik-Chervonenkis dimension* of  $\mathcal{H}$  as

$$VCD(\mathcal{H}) = \text{size of the largest shattered set} \quad (3.7)$$

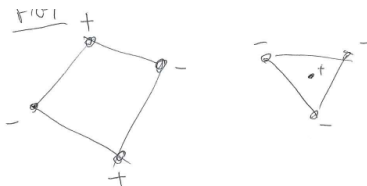
$$= \max_d \left\{ \exists S \subseteq X : |S| = d \quad \& \quad \Pi_{\mathcal{H}}(S) = \{0, 1\}^d \right\}. \quad (3.8)$$

**Example 3.14.** If  $\mathcal{H}$  is finite, then  $VCD(\mathcal{C}) \leq \ln|\mathcal{C}|$  because shattered  $d$  points requires  $2^d$  functions. Equivalently, if  $VCD(\mathcal{C}) = d$ , then we must have  $|\mathcal{C}| \geq 2^d$ .

**Example 3.15** (Rectangles in  $\mathbb{R}^2$ ). The VCD dimension is 4 in this case because the rectangular convex hull of a set of points is always given by the four extremal points in each direction. So there are always impossible labelings with 5 points. In general, the same kind of argument shows that the VCD dimension is  $2d$ .

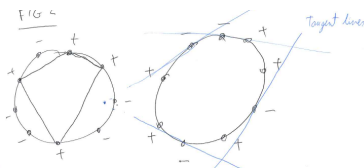
## 4 Lecture 4: 2017.02.13

**Example 4.1** (Hyperplanes in  $\mathbb{R}^2$ ). It is clear that the VCD is at least 3 because given three points, we can always choose a line that separates the points in any way we want. In fact, the VCD is exactly 3.



In general, for hyperplanes in  $\mathbb{R}^d$  the VCD dimension is  $d + 1$ .

**Example 4.2** (Convex n-gons in  $\mathbb{R}^2$ ). In this case,  $VCD = 2n + 1$ . For example, if  $n = 4$ ,  $VCD = 9$ .



*Remark 4.3.* Note that in these geometric situations, the VCD dimension is linearly related to the free parameters in our model. This is generally true (more or less). This supports the idea that we need more samples than parameters to be able to start learning something significant.

*Remark 4.4* (Monotone conjunctions over  $\{0, 1\}^m$  e.g.  $x_1 \wedge x_5 \wedge x_7$ ). We claim that  $VCD \geq n$ . Indeed, consider the bit vectors:  $\langle 01 \dots 1 \rangle, \langle 101 \dots 1 \rangle, \dots, \langle 1 \dots 10 \rangle$ . We can always find a conjunction with the desired label. Namely,  $x_1 \wedge \dots \wedge x_i^\vee \wedge \dots \wedge x_n$  or its negation as needed.

### Overloaded Notation

$$\Pi_{\mathcal{C}}(m) := \max_{S: |S|=m} \left\{ |\Pi_{\mathcal{C}}(S)| \right\} \quad (4.1)$$

$$= \max. \# \text{ of labelings induced by } \mathcal{C} \text{ on } m \text{ points.} \quad (4.2)$$

Note that we have the following:

- $\Pi_{\mathcal{C}}(m) \leq 2^m$ .
- If  $m \leq VCD(\mathcal{C})$ ,  $\Pi_{\mathcal{C}}(m) = 2^m$ .
- If  $m > VCD(\mathcal{C})$ ,  $\Pi_{\mathcal{C}}(m) < 2^m$ .

**Theorem 4.5.** Let  $VCD(\mathcal{C}) = d$ . Then  $\Pi_{\mathcal{C}}(m) = O(m^d)$ .

*Remark 4.6.* As the sample size becomes large compared to the dimension, the growth rate of  $\Pi_{\mathcal{C}}(m)$  is sub-exponential. In fact, the fraction of possible behaviors grows like  $\frac{1}{2^m} m^d$  which tends to 0 as  $m \rightarrow \infty$ .

### Proof of Theorem 4.5

Consider  $\phi_d(m) := \phi_d(m-1) + \phi_{d-1}(m-1)$  with boundary conditions  $\phi_0(m) = \phi_d(0) = 1$ . We first prove the following lemma.

**Lemma 4.7.** *If  $VCD(\mathcal{C}) = d$ , then for any  $m$  we have  $\Pi_{\mathcal{C}}(m) \leq \phi_d(m)$ .*

*Proof.* We use a double induction argument. The base cases are clear so assume the lemma is true for any  $d' \leq d$  and  $m' \leq m$  with at least one the inequalities being  $<$ . Consider the set  $S = \{x_1, \dots, x_m\}$ . Recall that  $\Pi_{\mathcal{C}}(S)$  is just the possible labelings of  $S$  that can be realized with concepts in  $\mathcal{C}$ . Now consider all the possible labelings that are possible for  $\{x_1, \dots, x_{m-1}\} = S \setminus \{x_m\}$ . For each of these, there are exactly 2 potential labelings for  $S$  depending on whether  $x_m = 0$  or  $x_m = 1$ . Then, by the induction hypothesis, we have

$$|\Pi_{\mathcal{C}}(S \setminus \{x_m\})| \leq \Pi_{\mathcal{C}}(m-1) \leq \phi_d(m-1). \quad (4.3)$$

Next consider the labelings of  $S \setminus \{x_m\}$  that appear twice in the labelings of  $S$ . That is, those such that both  $x_m = 0$  and  $x_m = 1$  are possible labelings. Then we can shatter at most  $d-1$  points with  $S \setminus \{x_m\}$  because otherwise we would violate  $VCD(\mathcal{C}) = d$ . Finally, by induction hypothesis we have that

$$\# \text{ of labelings of } S \setminus \{x_m\} \leq \phi_{d-1}(m-1). \quad (4.4)$$

Combining these two bound we obtain the result.  $\square$

**Claim 4.8.**  $\phi_d(m) = \sum_{i=0}^d \binom{m}{i}$ .

*Proof.* By induction,

$$\phi_d(m) = \phi_d(m-1) + \phi_{d-1}(m-1) = \sum_{i=0}^d \binom{m-1}{i} + \sum_{i=0}^{d-1} \binom{m-1}{i} \quad (4.5)$$

$$= \sum_{i=0}^d \binom{m-1}{i} + \sum_{i=0}^d \binom{m-1}{i-1} \quad (4.6)$$

$$= \sum_{i=0}^d \left[ \binom{m-1}{i} + \binom{m-1}{i-1} \right] = \sum_{i=0}^d \binom{m}{i}. \quad (4.7)$$

$\square$

Now the proof of Theorem 4.5 follows by noting that  $\sum_{i=0}^d \binom{m}{i} = O(m^d)$ .

## 4.1 Combinatorial $\rightarrow$ Probabilistic

Fix target  $c \in \mathcal{C}$  and a distribution  $D$ . Thinking of  $c$  as a set, we can consider  $\mathcal{H} \Delta c = \{h \Delta c : h \in \mathcal{H}\}$ , which are the *error regions*.

**Definition 4.9.** We say that  $S$  of  $m$  points is a  $\varepsilon$ -net for  $\mathcal{H}\Delta c$  if for all  $r \in \mathcal{H}\Delta c$  such that  $D[r] \geq \varepsilon$ , then we have  $S \cap r \neq \emptyset$ .

*Remark 4.10.*  $VCD(\mathcal{H}\Delta c) = VCD(\mathcal{H})$ . This is because taking the symmetric difference corresponds to a XOR operation.

Consider the auxiliary function on  $S = \{x_1, \dots, x_m\}$

$$A(S) := \begin{cases} 1 & \text{if } S \text{ is not an } \varepsilon\text{-net for } \mathcal{H}\Delta c \\ 0 & \text{else} \end{cases}, \quad (4.8)$$

$$B(S, S') := \begin{cases} 1 & \text{if } \exists r \in \mathcal{H}\Delta c \text{ st. } D[r] \geq \varepsilon, r \cap S = \emptyset, |r \cap S'| \geq \frac{\varepsilon m}{2} \\ 0 & \text{else} \end{cases}. \quad (4.9)$$

*Remark 4.11.* Event  $B(S, S')$  occurs if  $\exists r \in \mathcal{H}\Delta c$  st.  $r$  is “hit”:

- $l \geq \frac{\varepsilon m}{2}$  times in  $S \cup S'$  but *all*  $l$  fall in  $S'$ .

$$\mathbb{P}_{S, S'}[B(S, S')] = \mathbb{P}_{S, S'}[B(S, S') \mid A(S)] \cdot \mathbb{P}_{S, S'}[A(S)] \quad (4.10)$$

$$\geq \frac{1}{2} \mathbb{P}_{S, S'}[A(S)], \quad (4.11)$$

where we used Chebyshev’s inequality on the last line.

**Draw  $\langle S, S' \rangle$  of  $2m$  points in 2 steps:**

1. Draw  $2m$  points randomly from the distribution  $D$ .  $x_1, \dots, x_m, x_{m+1}, \dots, x_{2m}$   
- # possible labelings of the  $2m$  points by  $\mathcal{H}\Delta c$  fixed  $\leq \phi_d(2m)$ .
2. Split the  $2m$  points randomly into  $S$  and  $S'$ .

**Claim 4.12.**  $\mathbb{P}_{\text{permutation}}[B(S, S')] = \frac{\binom{m}{l}}{\binom{2m}{l}} \leq 2^{-l}$ .

*Proof.* Simple exercise in combinatorics. □

Then we have that

$$\mathbb{P}[B(S, S')] \leq \phi_d(2m) 2^{-\frac{\varepsilon m}{2}} \leq C_0(2m)^d 2^{-\frac{\varepsilon m}{2}} \leq e^{C_1[d \log(m) - \varepsilon m]}. \quad (4.12)$$

We have just proven the following theorem.

**Theorem 4.13.** As long as  $m \geq C_2 \left( \frac{1}{\varepsilon} \log \left( \frac{1}{\delta} \right) + \frac{d}{\varepsilon} \log \left( \frac{1}{\varepsilon} \right) \right)$ , where  $VCD(\mathcal{H}) = d$ , we have that with probability  $\geq 1 - \delta$ , any consistent  $h \in \mathcal{H}$  has error  $\leq \varepsilon$  with respect to  $c$  and  $D$ .

## 5 Lecture 5: 2017.02.20

### Outline for today:

1. VC theory recap
2. Matching lower bound
3. VC++:
  - unrealizable/agnostic case
  - ERM/SRM
  - Rademacher complexity, sparsity, etc.
  - general loss functions
4. Learning with noise

### 5.1 VC theory review

Recall from last time that we defined  $\Pi_{\mathcal{C}}(m)$  as the max. # of labelings induced by  $\mathcal{C}$  on  $m$  points, by definition we have the naive bound  $\Pi_{\mathcal{C}}(m) \leq 2^m$ . One of the highlights is that we were able to prove that

$$\Pi_{\mathcal{C}}(m) \leq \phi_d(m) \leq O(m^d). \quad (5.1)$$

This bound is saying that past the VCD dimension, the number of possible labelings we can achieve is only an increasingly small fraction of the total number ( $2^m$ ). This was the key step to prove the following theorem

**Theorem 5.1.** *Given  $c \in \mathcal{C}$  and a distribution,  $D$ , over  $X$ . If we take  $m \geq c_0 \frac{d}{\varepsilon} \ln \frac{1}{\delta}$  examples, where  $VCD(\mathcal{C}) = d$ , then, with probability  $\geq 1 - \delta$ , any  $h \in \mathcal{C}$  that is consistent with the sample has  $Err(h) \leq \varepsilon$ .*

*Remark 5.2.* This is a purely “information theory” statement, we are ignoring how hard it is to actually find an appropriate  $h$ .

### 5.2 Matching lower bound

The converse of the theorem above is not true, in other words, we can find distributions such that it is not necessary to take that many examples. For example, consider the case of the triangles in dimension 2 and a distribution that has support on a single point. This motivates the next result.

**Theorem 5.3.** *Given  $\mathcal{C}$  with  $VCD(\mathcal{C}) = d$ , there exists  $D$  a distribution over  $X$  such that  $\Omega(d/\varepsilon)$  examples are required to learn with error  $\leq \varepsilon$ .*

*Proof.* Consider  $\{x_1, \dots, x_d\}$  a shattered set and let  $D$  be the uniform distribution over this set. Now, since the set is shattered, we can choose a function  $c_i \in \mathcal{C}$  with  $i = 1, \dots, 2^d$  for each of the possible labeling of these  $d$  points. Now we randomly choose the target concept,  $c$ , among these  $c_i$ . This is equivalent to flipping a fair coin  $d$  times to determine the labeling induced by  $c$  on  $S$ .

Now we let  $L$  be any  $PAC$  learning algorithm with the data above. Given error parameter  $\varepsilon \leq \frac{1}{8}$  suppose we only draw  $m < d$  examples. Say that the number of distinct points that we have seen is  $m' \leq m$ , for the remaining of the points the problem is equivalent to predicting a fair coin toss. Therefore, the expected error on the whole sample is  $\frac{\frac{d-m'}{2}}{d} = \frac{d-m'}{2d}$ . By **Chebyshev's inequality**,

$$\mathbb{P}\left(\text{Error} \geq \frac{d-m'}{4d}\right) \leq \frac{1}{2}. \quad (5.2)$$

In particular, if  $m = \frac{d}{2}$ , the error is at least  $\frac{1}{8}$  with probability at least  $\frac{1}{2}$  and the algorithm fails the conclusion of the theorem when  $c$  is chosen randomly. This implies that there exist some target  $c$  on which  $L$  fails, so we need at least  $\Omega(d)$  sample complexity lower bound.

Finally, to incorporate the  $\varepsilon$  onto the lower bound, we need only scale the distribution above. Modify  $D$  so that the point  $x_1$  has probability  $1 - 8\varepsilon$  and let the rest have probability  $\frac{8\varepsilon}{d-1}$ . This results in needing to draw more points to be in the same set up as before. Details are left as exercise or can be checked on the textbook in page 63.  $\square$

## 5.3 VC++

### 5.3.1 Unrealizable/Agnostic setting

So far we have assumed that the target  $c \in \mathcal{C}$ . That is, that the truth is perfectly represented by our chosen class. Now we drop this assumption. We have a *hypothesis* class  $\mathcal{H}$  but  $c \notin \mathcal{H}$ , i.e. the target might not be representable by  $\mathcal{H}$ . The first thing we need to lower our expectations for learning.

**Definition 5.4.**  $h^* := \underset{h \in \mathcal{H}}{\operatorname{argmin}} \{ \text{Err}(h) \}$ , where  $\text{Err}(h) = \mathbb{P}_{X \sim D}[h(x) \neq c(x)]$ . We also define  $\varepsilon^* = \text{Err}(h^*) =$  best possible error in  $\mathcal{H}$ .

Let's set up some additional notation. Given  $h \in \mathcal{H}$  and a labeled sample

$$S = \{ \langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle \}, \quad (5.3)$$

we have some different types of errors:

- $\text{Err}(h)$ , the *true error* as above
- $\widehat{\text{Err}}(h) = \frac{1}{m} |\{ \langle x_i, y_i \rangle \in S : h(x_i) \neq y_i \}|$ , the *empirical error*.

**Theorem 5.5.** *Given any target  $c$  (not necessarily in  $\mathcal{H}$ ) and a distribution  $D$  over  $X$ , if we take  $m \geq c_0 \frac{d}{\varepsilon^2} \ln \frac{1}{\delta}$  where  $VCD(\mathcal{H}) = d$ , then with probability  $\geq 1 - \delta$ , for all  $h \in \mathcal{H}$  we have uniform convergence*

$$\left| \widehat{\text{Err}}(h) - \text{Err}(h) \right| \leq \varepsilon. \quad (5.4)$$

*Proof.* According to Prof. Kearns is a modification of the analogous result where  $c \in \mathcal{H}$ . He mentioned that some extra materials might be uploaded to the course website.  $\square$

*Remark 5.6.* To apply this theorem, in applications, we find/compute  $\hat{h}^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \left\{ \widehat{Err}(h) \right\}$ .

By the theorem, we have  $\left| \widehat{Err}(\hat{h}^*) - Err(\hat{h}^*) \right| \leq \varepsilon$  and  $\left| \widehat{Err}(h^*) - Err(h^*) \right| \leq \varepsilon$ . This implies that,

$$\left| Err(\hat{h}^*) - Err(h^*) \right| \leq 3\varepsilon. \quad (5.5)$$

### 5.3.2 Structural Risk Minimization

Fix some target  $c$ , a distribution  $D$  and a sample of size  $m$   $S$ . Suppose we have  $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots \subset \mathcal{H}_d \subset \dots$ . For example, each hypothesis could be a neural network with an increasing number of hidden layers. Note that the VCD dimension can only increase along the nested sequence. For simplicity, let's assume that  $VCD(\mathcal{H}_d) = d$ . Define

$$Err(d) := \min_{h \in \mathcal{H}_d} \{Err(h)\} \quad (5.6)$$

$$= \text{best true error in } \mathcal{H}_d. \quad (5.7)$$

$$\widehat{Err}(d) := \min_{h \in \mathcal{H}_d} \{\widehat{Err}(h)\} \quad (5.8)$$

$$= \text{best empirical error on } S \text{ in } \mathcal{H}_d. \quad (5.9)$$

*Remark 5.7.* From the theorem above,  $|Err(d) - \widehat{Err}(d)| \lesssim \sqrt{\frac{d}{m}}$

If we use minimization of the empirical error as the criterion to choose our model, we might run into trouble in the form of overfitting. Basically, the empirical error can only decrease as the complexity of the model increases but it can deviate from the true error. Luckily for us, the theorem above gives a way to correct for this by instead finding

$$\min_d \left\{ \widehat{Err}(d) + \sqrt{\frac{d}{m}} \right\}. \quad (5.10)$$

### 5.3.3 Refinements

In the theorem from last time,  $\Pi_C(S) \rightarrow \Pi_C(m) \leq O(m^d)$ . But if one is more careful in the proof, it is possible to arrive to a bound on  $\mathbb{E}_{S \sim D, c} [\ln |\Pi_C(S)|]$  where  $|S| = m$ , the *VC entropy*. (cf. Book)

## 5.4 General loss functions

The assumptions are as follows:

- observations  $z \sim D$  iid. E.g.  $z = \langle x, y \rangle$  with  $y \in \{0, 1\}$ ;  $z = \langle x, y \rangle$  with  $y \in \mathbb{R}$ ;  $z = x$ .
- models  $h \in \mathcal{H}$ .

- $\mathbb{R}$ -valued loss function,  $L(h, z)$ . E.g.
  1. Classification:  $L(h, \langle x, y \rangle) = 1$  if  $h(x) \neq y$ , or 0 otherwise.
  2. Linear regression:  $L(h, \langle x, y \rangle) = (h(x) - y)^2$
  3.  $L(h, x) = \ln \frac{1}{h(x)}$

We'd like to minimize  $\mathbb{E}_{z \sim D}[L(h, z)]$

*Remark 5.8.* Although we are not going to touch on that, there is a theorem like the one we proved for this setting as well.

#### 5.4.1 Analog to VCD

Take  $d$  observations  $z_1, \dots, z_d$  and consider the following set  $T = \{\langle L(h, z_1), \dots, L(h, z_d) \rangle : h \in \mathcal{H}\}$ . The analog to the VC dimension would be to look at this cloud of points and see if it is “space-filling” or that not all points lie on some lower dimensional subspace. The criterion is that you want all the  $d$ -dimensional orthants to be intersected by  $T$ .