

Understanding Machine Learning: Exercises

Somnath Sikdar

August 9, 2019

Contents

	About this Document	2
3	A Formal Learning Model	3
4	Learning via Uniform Convergence	9
5	The No-Free-Lunch Theorem	11
6	The VC-Dimension	13

What this is About

These notes are my attempt to understand and work out material from the textbook *Understanding Machine Learning* by Shai Shalev-Shwartz and Shai Ben-David.

Chapter 3

A Formal Learning Model

Notes on Chapter 3

The main concept introduced here is that of agnostic PAC learnability. It helps to review the definitions of both PAC learnability with the realizability assumption and that of agnostic PAC learnability.

Definition 1 (PAC Learnability). Fix a domain \mathcal{X} , a range \mathcal{Y} and let \mathcal{H} be a set of functions from $\mathcal{X} \rightarrow \mathcal{Y}$. The class \mathcal{H} is PAC learnable if there exists a function $m_{\mathcal{H}}: (0, 1) \times (0, 1) \rightarrow \mathbf{N}$ and a learning algorithm \mathcal{A} such that the following holds: for all $\epsilon, \delta \in (0, 1)$, all labeling functions $f: \mathcal{X} \rightarrow \mathcal{Y}$ and all distributions \mathcal{D} over \mathcal{X} such that \mathcal{H} is realizable wrt \mathcal{D} and f , if \mathcal{A} is presented with a sample of at least $m_{\mathcal{H}}(\epsilon, \delta)$ examples drawn iid from \mathcal{D} , \mathcal{A} returns a hypothesis h_S such that

$$\Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D},f}(h_S) \leq \epsilon\} > 1 - \delta.$$

In the definition above, the sample complexity function $m_{\mathcal{H}}$ must work for all possible labeling functions f and distributions \mathcal{D} as long as \mathcal{H} is realizable (wrt this labeling function and distribution). A learning task is completely specified by $(\mathcal{X}, \mathcal{Y}, f, \mathcal{D})$. Intuitively, what this definition states is that a hypothesis class is PAC learnable if there exists a learner which when given a sufficiently large number of training examples can approximate the true labeling function f with high probability for all learning tasks on the domain $\mathcal{X} \times \mathcal{Y}$ that satisfy the realizability condition.

Definition 2 (Agnostic PAC Learnability). A class \mathcal{H} is PAC learnable if there exists a function $m_{\mathcal{H}}: (0, 1) \times (0, 1) \rightarrow \mathbf{N}$ and a learning algorithm \mathcal{A} such that the following holds: for all $\epsilon, \delta \in (0, 1)$ and all distributions \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ such that, if \mathcal{A} is presented with a sample of at least $m_{\mathcal{H}}(\epsilon, \delta)$ examples drawn iid from \mathcal{D} , \mathcal{A} returns a hypothesis h_S such that

$$\Pr_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(h_S) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon \right\} > 1 - \delta.$$

In the agnostic setting, a learning task is completely specified by the triplet $(\mathcal{X}, \mathcal{Y}, \mathcal{D})$ and a hypothesis class is agnostic PAC learnable if there exists a learner which when given a sufficiently many training examples outputs a “good enough” hypothesis with high probability.

The goodness of the hypothesis is measured with respect to the best hypothesis of the class \mathcal{H} .

Agnostic PAC learnability presents a *stronger* requirement as one must be able to learn *any* distribution as distinct from PAC learnability where one must be able to learn distributions for which the hypothesis class is realizable. As such, agnostic PAC learnability implies PAC learnability.

Exercise 3.1

Let $m_{\mathcal{H}}(\epsilon, \delta)$ be the sample complexity of a PAC-learnable hypothesis class \mathcal{H} for a binary classification task. For a fixed δ , let $0 < \epsilon_1 \leq \epsilon_2 < 1$ and suppose that $m_{\mathcal{H}}(\epsilon_1, \delta) < m_{\mathcal{H}}(\epsilon_2, \delta)$. Then when running the learning algorithm on $m_{\mathcal{H}}(\epsilon_1, \delta)$ i.i.d examples, we obtain a hypothesis h , which with probability at least $1 - \delta$ has a true error $L_{\mathcal{D},f}(h) \leq \epsilon_1 \leq \epsilon_2$. This implies that for the (ϵ_2, δ) combination of parameters, we can bound the true error of h by ϵ_2 by using a smaller number of i.i.d examples than $m_{\mathcal{H}}(\epsilon_2, \delta)$. This contradicts the minimality of the sample complexity function. Hence we must have $m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$.

Next suppose that $0 < \delta_1 \leq \delta_2 < 1$ and that $m_{\mathcal{H}}(\epsilon, \delta_1) < m_{\mathcal{H}}(\epsilon, \delta_2)$, where ϵ is fixed in advance. Then with $m_{\mathcal{H}}(\epsilon, \delta_1)$ i.i.d examples, the learner outputs a hypothesis h which with probability at least $1 - \delta_1 \geq 1 - \delta_2$ has a true error of at most ϵ . This implies that for the (ϵ, δ_2) combination of parameters, we can bound the true error of h by ϵ by using a smaller number of i.i.d examples than $m_{\mathcal{H}}(\epsilon, \delta_2)$. This again contradicts the minimality of the sample complexity function. Hence we must have $m_{\mathcal{H}}(\epsilon, \delta_1) \geq m_{\mathcal{H}}(\epsilon, \delta_2)$.

Exercise 3.2

Given a sample S , we output a hypothesis h_S with the property that $\forall x \in S_x$,

$$h_S(x) = \begin{cases} 1, & \text{if } (x, 1) \in S \\ 0, & \text{otherwise} \end{cases}$$

For any sample S , this hypothesis has an empirical loss of 0. Note that h_S disagrees with the true labeling function f in at most one point $z \in \mathcal{X}$. It's true loss is therefore $\Pr_{x \sim \mathcal{D}}\{f(x) \neq h_S(x)\} = \Pr_{\mathcal{D}}\{z\} := p_z$.

The true loss of h_S will be 0 if $(z, 1) \in S$. Therefore the probability of getting a “bad” sample is $\Pr_{S \sim \mathcal{D}^m}\{(z, 1) \notin S\}$. Let $z^* \in \mathcal{X}$ be a point at which $(1 - p_z)^m$ is maximized. Since $(1 - p_{z^*})^m \leq e^{-mp_{z^*}}$ and since we want the probability of picking a bad sample to be at most δ , we want $e^{-mp_{z^*}} < \delta$, which gives us the sample size to be:

$$m > \frac{\log(1/\delta)}{p_{z^*}} \quad (3.1)$$

Depending on the value of the error bound ϵ , there are two situations to consider. If $\epsilon \geq p_{z^*}$, then even a sample of size one will guarantee that the true error of h_s is at most ϵ .

However if $\epsilon < p_{z^*}$ then we can then use this in (3.1) to obtain:

$$m > \frac{\log(1/\delta)}{\epsilon}.$$

Thus the sample complexity is $m_{\mathcal{H}}(\epsilon, \delta) = \max \left\{ 1, \frac{\log(1/\delta)}{\epsilon} \right\}$.

Exercise 3.3

Here $\mathcal{X} = \mathbf{R}^2$ and $\mathcal{Y} = \{0, 1\}$. The hypothesis class \mathcal{H} is the set of concentric circles in \mathbf{R}^2 centered at the origin. Assuming realizability, this implies that the true labeling function $f = h_r$ for some $r \in \mathbf{R}_+$. Thus f assigns the label 1 to any point (x, y) that is within a distance of r from the origin and 0 otherwise.

Given any sample S , let $q \in \mathbf{R}_+$ be the minimum real number such that all $(x, y) \in S_x$ with a label of 1 are included in a circle centered at the origin with radius q . The output of the ERM procedure is h_q . The empirical error of h_q is zero, but it's true error is:

$$\Pr_{(x,y) \sim \mathcal{D}} \{(x, y) \in C_r \setminus C_q\}$$

where C_r and C_q are concentric circles centered at the origin with radius r and q respectively. Given an $\epsilon > 0$, let $t \in \mathbf{R}_+$ be such that

$$\epsilon = \Pr_{(x,y) \sim \mathcal{D}} \{(x, y) \in C_r \setminus C_t\}.$$

That is, we choose t so that the true error matches the probability of picking anything inside the ring described by the circles C_r and C_t . Then the probability that we fail to choose any point in this ring in an i.i.d sample of size m is $(1 - \epsilon)^m \leq e^{-\epsilon m}$. This is the probability that we are handed a “bad” sample. Upper bounding this by δ , we obtain that $m > \log(1/\delta)/\epsilon$.

Now a sample of size at least $\log(1/\delta)/\epsilon$ has with probability at least $1 - \delta$ a point from $C_r \setminus C_t$, and hence the true error of the resulting ERM hypothesis is at most ϵ . Hence the sample complexity is upper bounded by $\lceil \log(1/\delta)/\epsilon \rceil$.

Exercise 3.4

In this example, $\mathcal{X} = \{0, 1\}^d$, $\mathcal{Y} = \{0, 1\}$ and the hypothesis class \mathcal{H} is the set of all conjunctions over d Boolean variables. Since there are $\sum_{i=0}^d \binom{d}{i} 2^i = 3^d$ Boolean conjunctions over d Boolean variables, the hypothesis class is finite. Hence the sample complexity is

$$\begin{aligned} m_{\mathcal{H}}(\epsilon, \delta) &= \left\lceil \frac{\log(\mathcal{H}/\delta)}{\epsilon} \right\rceil \\ &= \left\lceil \frac{d \cdot \log 3 + \log(1/\delta)}{\epsilon} \right\rceil \end{aligned}$$

To prove that the class \mathcal{H} is PAC-learnable, it suffices to exhibit a polynomial-time algorithm that implements the ERM rule. The algorithm outlined in Figure 3.1 starts with

```

procedure PACBOOLEAN( $S$ )  $\triangleright$   $S$  is the sample set with elements  $\langle (a_1, \dots, a_d), b \rangle$ , where
 $(a_1, \dots, a_d) \in \{0, 1\}^d$  and  $b \in \{0, 1\}$ 
   $f \leftarrow x_1 \wedge \bar{x}_1 \wedge \dots \wedge x_d \wedge \bar{x}_d$ 
  for each  $\langle (a_1, \dots, a_d), b \rangle \in S$  with  $b = 1$  do
    for  $j$  in  $[1, \dots, d]$  do
      if  $a_j = 1$  then
        Delete  $\bar{x}_j$  from  $f$ , if it exists in the formula
      else
        Delete  $x_j$  from  $f$ , if it exists in the formula
      end if
    end for
  end for
  return  $f$ 
end procedure

```

Figure 3.1: Learning Boolean conjunctions

the formula $x_1 \wedge \bar{x}_1 \wedge \dots \wedge x_d \wedge \bar{x}_d$. It runs through the positive examples in the sample S and for each such example, it adjusts the formula so that it satisfies the assignment given in the example. At the end of this procedure, the modified formula satisfies all positive examples of S . The time taken is $O(d \cdot |S|)$.

What may not be immediately apparent is that the formula returned by the algorithm satisfies all negative examples too. This is clear when the sample S has *no* positive examples to begin with as every assignment to $x_1 \wedge \bar{x}_1 \wedge \dots \wedge x_d \wedge \bar{x}_d$ results in a 0. The point is that if there is even *one* positive example, for each $1 \leq i \leq d$, the algorithm eliminates either x_i or \bar{x}_i depending on the assignment. That is, it eliminates half of the literals on seeing that one example and the modified formula f contains the literals of the true labeling function along with possibly others. Now the literals of the true labeling function produce a 0 on all negative examples and so does f . Hence the sampling error of the function returned by the algorithm is 0.

Exercise 3.5

The first thing to verify is that $\bar{\mathcal{D}}_m$ is a distribution. This is easy since for all $x \in \mathcal{X}$, $\bar{\mathcal{D}}_m(x) \geq 0$ and

$$\begin{aligned}
 \int_{x \in \mathcal{X}} \bar{\mathcal{D}}_m(x) dx &= \frac{1}{m} \sum_{i=1}^m \int_{x \in \mathcal{X}} \mathcal{D}_i(x) dx \\
 &= \frac{1}{m} \sum_{i=1}^m 1 \\
 &= 1.
 \end{aligned}$$

Fix an accuracy parameter $\epsilon > 0$. As in the text, define the set of bad hypotheses to

be $\mathcal{H}_B = \{h \in \mathcal{H} : L_{\bar{\mathcal{D}}_{m,f}}(h) > \epsilon\}$ and let $\mathcal{M} = \{S|_x : \exists h \in \mathcal{H}_B, L_S(h) = 0\}$ be the set of misleading samples. Since we assume realizability, any hypothesis h output by the ERM procedure has $L_S(h) = 0$. Thus the event $L_{\bar{\mathcal{D}}_{m,f}}(h) > \epsilon$ and $L_S(h) = 0$ happens only when $S|_x \in \mathcal{M}$. Hence,

$$\begin{aligned}
\Pr_{\forall i: x_i \sim \mathcal{D}_i} \{S|_x \in \mathcal{M}\} &= \Pr_{\forall i: x_i \sim \mathcal{D}_i} \left\{ \bigcup_{h \in \mathcal{H}_B} \{S|_x : L_S(h) = 0\} \right\} \\
&\leq \sum_{h \in \mathcal{H}_B} \Pr_{\forall i: x_i \sim \mathcal{D}_i} \{S|_x : L_S(h) = 0\} \\
&= \sum_{h \in \mathcal{H}_B} \prod_{i=1}^m \Pr_{x_i \sim \mathcal{D}_i} \{f(x_i) = h(x_i)\} \\
&= \sum_{h \in \mathcal{H}_B} \prod_{i=1}^m (1 - L_{\mathcal{D}_i,f}(h)) \\
&\leq \sum_{h \in \mathcal{H}_B} \left[\frac{1}{m} \sum_{i=1}^m (1 - L_{\mathcal{D}_i,f}(h)) \right]^m \\
&\leq \sum_{h \in \mathcal{H}_B} [1 - L_{\bar{\mathcal{D}}_{m,f}}(h)]^m
\end{aligned}$$

The second-last inequality follows from the fact that the arithmetic mean of a set of numbers is at most their geometric mean. The quantity $\sum_{h \in \mathcal{H}_B} [1 - L_{\bar{\mathcal{D}}_{m,f}}(h)]^m$ is at most $|\mathcal{H}| \cdot (1 - \epsilon)^m$ which is at most $|\mathcal{H}| \cdot e^{-\epsilon m}$.

Exercise 3.6

Agnostic PAC-learnability implies PAC-learnability. Let \mathcal{H} be a set of functions from \mathcal{X} to $\{0, 1\}$ which is agnostic PAC-learnable wrt $\mathcal{X} \times \{0, 1\}$ and the 0-1 loss function with sample complexity $m_{\mathcal{H}}$. Let f be a labeling function and let $\mathcal{D}_{\mathcal{X}}$ be a distribution over \mathcal{X} for which the realizability assumption holds, that is, there exists $h \in \mathcal{H}$ such that $L_{\mathcal{D}_{\mathcal{X}},f}(h) = 0$.

Define a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ as follows: for all $x \in \mathcal{X}$, $\mathcal{D}((x, f(x))) = \mathcal{D}_{\mathcal{X}}(x)$ and $\mathcal{D}((x, 1 - f(x))) = 0$. Fix $\epsilon, \delta > 0$. Since \mathcal{H} is agnostic PAC-learnable, there exists a learner A which given a sample S of $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ iid examples generated by \mathcal{D} returns a hypothesis h_S such that

$$\Pr_{S \sim \mathcal{D}^m} \left\{ L_{\mathcal{D}}(h_S) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon \right\} > 1 - \delta.$$

Note that for any $h' \in \mathcal{H}$, we may write the loss $L_{\mathcal{D}}(h')$ as follows.

$$\begin{aligned}
L_{\mathcal{D}}(h') &= \Pr_{(x,y) \in \mathcal{D}} \{h'(x) \neq y\} \\
&= \Pr_{(x,y) \in \mathcal{D}} \{h'(x) \neq f(x)\} \\
&= L_{\mathcal{D}_{\mathcal{X}},f}(h').
\end{aligned}$$

The second equality above follows since the only points $(x, y) \in \mathcal{X} \times \{0, 1\}$ for which \mathcal{D} places a non-zero probability mass are those for which $y = f(x)$. Since we assume realizability, $\min_{h' \in \mathcal{H}} L_{\mathcal{D}_{\mathcal{X}}, f}(h') = 0$. Hence the hypothesis h_S returned by the learner A satisfies:

$$\Pr_{S|\mathcal{X} \sim \mathcal{D}_{\mathcal{X}}^m} \{L_{\mathcal{D}_{\mathcal{X}}, f}(h_S) \leq \epsilon\} > 1 - \delta,$$

which is the condition for successful PAC-learnability.

Exercise 3.7

Let us fix some notation. We assume that X and Y are random variables defined over the domains \mathcal{X} and $\{0, 1\}$, respectively. Let $\mathcal{D}_{X,Y}$ be a distribution over $\mathcal{X} \times \{0, 1\}$; let $\mathcal{D}_{Y|X}$, the conditional distribution of Y given X ; let \mathcal{D}_X be the marginal distribution of X over \mathcal{X} ; and, finally, let $\eta(x) = \Pr_{\mathcal{D}_{Y|X}} \{Y = 1 \mid X = x\}$.

The Bayes optimal classifier $f_{\mathcal{D}}$ may be written as:

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \eta(x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

Given any classifier $g: \mathcal{X} \rightarrow \{0, 1\}$, the risk of this classifier is

$$L_{\mathcal{D}}(g) = \Pr_{\mathcal{D}_{X,Y}} \{g(X) \neq Y\} = \int_{x \in \mathcal{X}} \Pr_{\mathcal{D}_{Y|X}} \{g(x) \neq Y \mid X = x\} \cdot \Pr_{\mathcal{D}_X} \{X = x\} dx. \quad (3.2)$$

We may write the first term of this integrand as follows (where all probabilities are with respect to the conditional distribution $\mathcal{D}_{Y|X}$):

$$\begin{aligned} \Pr \{g(x) \neq Y \mid X = x\} &= 1 - \Pr \{g(x) = Y \mid X = x\} \\ &= 1 - [\Pr \{g(x) = 1, Y = 1 \mid X = x\} + \Pr \{g(x) = 0, Y = 0 \mid X = x\}] \\ &= 1 - [\mathbf{1}_{g(x)=1} \cdot \Pr \{Y = 1 \mid X = x\} + \mathbf{1}_{g(x)=0} \cdot \Pr \{Y = 0 \mid X = x\}] \\ &= 1 - [\mathbf{1}_{g(x)=1} \cdot \eta(x) + \mathbf{1}_{g(x)=0} \cdot (1 - \eta(x))] \end{aligned}$$

Consider the difference $\Pr \{g(x) \neq Y \mid X = x\} - \Pr \{f_{\mathcal{D}}(x) \neq Y \mid X = x\}$. This may be written as:

$$\begin{aligned} &[\mathbf{1}_{f_{\mathcal{D}}(x)=1} \cdot \eta(x) + \mathbf{1}_{f_{\mathcal{D}}(x)=0} \cdot (1 - \eta(x))] - [\mathbf{1}_{g(x)=1} \cdot \eta(x) + \mathbf{1}_{g(x)=0} \cdot (1 - \eta(x))] \\ &= (\mathbf{1}_{f_{\mathcal{D}}(x)=1} - \mathbf{1}_{g(x)=1}) \cdot \eta(x) + (\mathbf{1}_{f_{\mathcal{D}}(x)=0} - \mathbf{1}_{g(x)=0}) \cdot (1 - \eta(x)). \end{aligned}$$

This last expression may be written as:

$$(\mathbf{1}_{f_{\mathcal{D}}(x)=1} - \mathbf{1}_{g(x)=1}) \cdot \eta(x) + (\mathbf{1}_{g(x)=1} - \mathbf{1}_{f_{\mathcal{D}}(x)=1}) \cdot (1 - \eta(x)).$$

Rearranging terms allows us to write this as:

$$2 \cdot (\mathbf{1}_{f_{\mathcal{D}}(x)=1} - \mathbf{1}_{g(x)=1}) \cdot \eta(x) + (\mathbf{1}_{g(x)=1} - \mathbf{1}_{f_{\mathcal{D}}(x)=1}). \quad (3.3)$$

We claim that this last expression is always non-negative. If $f_{\mathcal{D}}(x) = 0$ then $\eta(x) < 1/2$ and the above expression is non-negative. If $f_{\mathcal{D}}(x) = 1$ then $\eta(x) \geq 1/2$ and, in this case too, the expression is non-negative. The result follows by plugging in the difference $\Pr \{g(x) \neq Y \mid X = x\} - \Pr \{f_{\mathcal{D}}(x) \neq Y \mid X = x\}$ in the integral in (3.2).

Chapter 4

Learning via Uniform Convergence

Notes on Chapter 4

Given any hypothesis class \mathcal{H} and a domain $Z = \mathcal{X} \times Y$, let l be a loss function from $\mathcal{H} \times Z \rightarrow \mathbf{R}_+$. Finally let \mathcal{D} be a distribution over the domain Z . The risk of a hypothesis $h \in \mathcal{H}$ is

$$L_{\mathcal{D}}(h) = \Pr_{z \sim \mathcal{D}} \{l(h, z)\}$$

A training set S is ϵ -representative w.r.t Z , \mathcal{H} , Z and l if for all $h \in \mathcal{H}$, $|L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$. Thus any hypothesis on an ϵ -representative training set has an in-sample error that is close to their true risk.

If S is ϵ -representative, then the $\text{ERM}_{\mathcal{H}}(S)$ learning rule is guaranteed to return a good hypothesis. More specifically,

Lemma 1. *Fix a hypothesis class \mathcal{H} , a domain $Z = \mathcal{X} \times Y$, a loss function $l: \mathcal{H} \times Z \rightarrow \mathbf{R}_+$ and a distribution \mathcal{D} over the domain Z . Let S be an $\epsilon/2$ -representative sample. Then any output h_S of $\text{ERM}_{\mathcal{H}}(S)$ satisfies*

$$L_{\mathcal{D}}(h_S) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon$$

Therefore in order for the ERM rule to be an agnostic PAC-learner, all we need to do is to ensure that with probability of at least $1 - \delta$ over random choices of the training set, we end up with an $\epsilon/2$ -representative training sample. This requirement is baked into the definition of *uniform convergence*.

Definition 3. A hypothesis class \mathcal{H} is uniformly convergent wrt a domain Z and a loss function l , if there exists a function $m_{\mathcal{H}}^{\text{UC}}: (0, 1) \times (0, 1) \rightarrow \mathbf{N}$ such that for all $\epsilon, \delta \in (0, 1)$ and all distributions \mathcal{D} on Z , if a sample of at least $m_{\mathcal{H}}^{\text{UC}}(\epsilon, \delta)$ examples is chosen i.i.d from \mathcal{D} , then with probability $1 - \delta$, the sample is ϵ -representative.

By Lemma (1), if \mathcal{H} is uniformly convergent with function $m_{\mathcal{H}}^{\text{UC}}$, then it is agnostically PAC-learnable with sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta)$. In this case, the ERM paradigm is a successful agnostic PAC-learner for \mathcal{H} .

The other main result is that all finite hypothesis classes are uniformly convergent and hence agnostic PAC learnable.

Exercise 4.1

We first show that (1) \Rightarrow (2). For each $n \in \mathbf{N}$, define $\epsilon_n = 1/2^n$ and $\delta_n = 1/2^n$. Then by (1), for each $n \in \mathbf{N}$, there exists $m(\epsilon_n, \delta_n)$ such that $\forall m \geq m(\epsilon_n, \delta_n)$,

$$\Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) > \epsilon_n\} < \delta_n.$$

We can then upper bound $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]$ as follows:

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] &\leq \epsilon_n \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) \leq \epsilon_n\} + (1 - \epsilon_n) \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) > \epsilon_n\} \\ &\leq \epsilon_n \cdot (1 - \delta_n) + (1 - \epsilon_n) \cdot \delta_n \\ &\leq \frac{1}{2^{n-1}} - \frac{1}{2^{2n-1}}. \end{aligned}$$

The first inequality follows from the fact that the loss function is from $\mathcal{H} \times Z \rightarrow [0, 1]$, which allows us to upper bound the value of the error when $L_{\mathcal{D}}(h_S) > \epsilon_n$ by $1 - \epsilon_n$. As $n \rightarrow \infty$, $m \rightarrow \infty$ and $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] \rightarrow 0$, proving that (2) follows.

We next show that (2) \Rightarrow (1). Fix $\epsilon, \delta > 0$. Define $\delta' = \epsilon \cdot \delta$. Since

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] = 0,$$

there exists $m_1(\delta')$ such that for all $m \geq m_1(\delta')$ we have $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] < \delta'$. We now lower bound $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]$ as follows:

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] &= \int_0^1 x \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) = x\} dx \\ &\geq \int_{\epsilon}^1 x \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) = x\} dx \\ &\geq \epsilon \cdot \int_{\epsilon}^1 \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) = x\} dx \\ &= \epsilon \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) \geq \epsilon\}. \end{aligned}$$

Choose $m(\epsilon, \delta) := m_1(\epsilon \cdot \delta)$. Then for all $m \geq m(\epsilon, \delta)$, we have that $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] < \epsilon \cdot \delta$, from which it follows that:

$$\epsilon \cdot \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(h_S) \geq \epsilon\} < \epsilon \cdot \delta.$$

Condition (1) follows from this.

Chapter 5

The No-Free-Lunch Theorem

Notes on Chapter 5

Consider a binary classification task on a domain \mathcal{X} . Assume for the time being that \mathcal{X} is finite. In this case, the set \mathcal{H} of all functions from $\mathcal{X} \rightarrow \{0, 1\}$ is finite and is hence PAC-learnable with sample complexity $\leq \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$. Since $|\mathcal{H}| = 2^{|\mathcal{X}|}$, the sample complexity is $\frac{|\mathcal{X}| + \log(1/\delta)}{\epsilon} = O(|\mathcal{X}|)$.

The first question is what happens wrt PAC-learnability in this situation when we restrict the sample size? The No-Free-Lunch theorem shows that there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ and a labelling function $f: \mathcal{X} \rightarrow \{0, 1\}$ that learners who are constrained to use at most $|\mathcal{X}|/2$ training examples “cannot learn.” There is another way to interpret the No-Free-Lunch theorem: if the domain \mathcal{X} is *infinite*, then the set of all functions from \mathcal{X} to $\{0, 1\}$ is not PAC-learnable no matter what the sample size.

Thus the No-Free-Lunch theorem has two interpretations, first, as a lower bound result on the sample complexity of PAC-learning and, second, as the inability to PAC-learn arbitrary hypothesis classes.

Theorem 1. *Consider the task of binary classification over the domain \mathcal{X} wrt the 0-1 loss function. Let A be a learning algorithm that is constrained to use at most $m \leq |\mathcal{X}|/2$ training examples. Then there exists a distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$ and a function $f: \mathcal{X} \rightarrow \{0, 1\}$ such that*

1. $L_{\mathcal{D}}(f) = 0$
2. *with probability of at least $1/7$ over the choice of training examples chosen iid from \mathcal{D}^m , we have that $L_{\mathcal{D}}(A(S)) \geq 1/8$.*

Exercise 5.1

As the hint in the exercise suggests, let θ be a random variable that takes on values in the range $[0, 1]$ with expectation $E[\theta] \geq 1/4$. We want to show that $\Pr\{\theta \geq 1/8\} \geq 1/7$.

We start with Markov's inequality: for any nonnegative random variable X and $a > 0$,

$$\Pr \{X \geq a\} \leq \frac{E[X]}{a}.$$

This doesn't quite work when we substitute $\theta = X$ and $a = 1/8$. The trick here lies in observing that θ is bounded from above by 1, and hence, if we define $\xi = 1 - \theta$ then ξ is nonnegative and we can use Markov's inequality on ξ . Note that $E[\xi] = 1 - E[\theta]$, and hence by Markov's inequality,

$$\begin{aligned} \Pr \{\xi \geq a\} &\leq \frac{E[\xi]}{a} \\ 1 - \Pr \{\xi \geq a\} &\geq 1 - \frac{E[\xi]}{a} \\ \Pr \{\xi < a\} &\geq 1 - \frac{1 - E[\theta]}{a} \\ \Pr \{1 - \theta < a\} &\geq \frac{a - 1}{a} + \frac{E[\theta]}{a} \end{aligned}$$

At this point, we use the fact that $E[\theta] \geq 1/4$ to obtain: $\Pr \{\theta > 1 - a\} \geq \frac{a-1}{a} + \frac{1}{4a}$. Now if we substitute $1 - a = 1/8$, or $a = 7/8$, then we obtain:

$$\Pr \{\theta > 1/8\} \geq 1/7.$$

Exercise 5.2

The first algorithm, the one that picks only blood pressure and the BMI as features, is simpler in the sense that the hypothesis class to be learned is simpler. We would expect that this algorithm has a higher bias but a lower variance when compared to the second algorithm which is more feature rich. The second algorithm would probably explain the conditions of a heart attack better as it includes relevant features such as age and the level of physical activity into account. We would expect the second algorithm to have a lower bias but a higher variance because there may be a tendency to overfit on any given sample.

Since the sample complexity is higher for a more complicated hypothesis class, if the sample size is "small," then we might want to choose the first algorithm. If sample size is not a problem, then the second algorithm is probably better.

Chapter 6

The VC-Dimension

Notes on Chapter 6

We know that finite hypothesis classes are agnostic PAC learnable (and hence PAC learnable). What about infinite hypothesis classes? The first example is that of an infinite hypothesis class that is PAC learnable.

Example 1 (Threshold Functions). Let $\mathcal{X} = [0, 1]$ and $\mathcal{Y} = \{0, 1\}$. For $r \in [0, 1]$, define $h_r: \mathcal{X} \rightarrow \mathcal{Y}$ as:

$$h_r(x) = \begin{cases} 0 & \text{if } x \leq r \\ 1 & \text{if } x > r \end{cases}$$

Let \mathcal{H}_{thr} be the set of all threshold functions h_r for $r \in [0, 1]$. Since \mathcal{H}_{thr} is not finite, it is not immediately obvious whether it is PAC learnable (in the realizable case).

Fix $\epsilon, \delta \in (0, 1)$. Let $f = h_s$ be the true labeling function where $s \in [0, 1]$ and let \mathcal{D} be the underlying distribution over the domain $[0, 1]$. Let $s_0 \in [0, s)$ and $s_1 \in [s, 1]$ be numbers such that

$$\mathcal{D}\{x \in [s_0, s)\} = \epsilon = \mathcal{D}\{x \in [s, s_1]\}$$

If $\mathcal{D}\{[0, s)\} < \epsilon$, then set $s_0 = 0$; similarly, if $\mathcal{D}\{[s, 1]\} < \epsilon$, set $s_1 = 1$. Since \mathcal{D} is a distribution, it must place a probability mass of ϵ either to the left or to the right of s .

Given a sample S , let $t_0 = \max\{t: (t, 0) \in S\}$ and $t_1 = \min\{t: (t, 1) \in S\}$. The ERM algorithm outputs h_p , where $p \in (t_0, t_1)$. In particular, if the sample presented to the ERM algorithm is such that $s_0 \leq t_0$ and $t_1 \leq s_1$, then hypothesis h_p returned by the ERM algorithm will incur a loss of $L_{\mathcal{D}}(h_p) \leq \epsilon$.

Thus the probability that the hypothesis $\text{ERM}(S)$ output by the ERM algorithm has a loss greater than ϵ on a sample S of size m is:

$$\begin{aligned} \Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(\text{ERM}(S)) > \epsilon\} &= \Pr_{S \sim \mathcal{D}^m} \{S: t_0 < s_0 \vee s_1 < t_1\} \\ &\leq \Pr_{S \sim \mathcal{D}^m} \{S: S|_x \cap [s_0, s) = \emptyset\} + \Pr_{S \sim \mathcal{D}^m} \{S: S|_x \cap [s, s_1] = \emptyset\} \\ &\leq 2 \cdot (1 - \epsilon)^m \\ &\leq 2 \cdot e^{-\epsilon m} \end{aligned}$$

Setting the last expression to be at most δ , we obtain that $m > \frac{1}{\epsilon} \cdot \log \frac{2}{\delta}$. Hence if we have samples of size at least $\frac{1}{\epsilon} \cdot \log \frac{2}{\delta}$,

$$\Pr_{S \sim \mathcal{D}^m} \{L_{\mathcal{D}}(\text{ERM}(S)) \leq \epsilon\} \geq 1 - \delta,$$

which is the condition for PAC learnability.

The second example shows that there are infinite hypothesis classes that are not PAC learnable at least by using an ERM strategy.

Example 2 (Identity Function for Finite Sets). Let $\mathcal{X} = \mathbf{R}$ and $\mathcal{Y} = \{0, 1\}$. Given a set $A \subseteq \mathcal{X}$, define h_A as follows:

$$h_A = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise} \end{cases}$$

Let $\mathcal{H}_{\text{finite}}$ be the set of all such functions h_A for *finite* subsets A of \mathbf{R} along with the function h_1 which maps every point in \mathbf{R} to 1. We claim that $\mathcal{H}_{\text{finite}}$ is not PAC learnable by an ERM algorithm.

Consider the case when the true labeling function $f = h_1$, the all-ones function on \mathbf{R} and \mathcal{D} is the uniform distribution on $[0, 1]$. Since $f \in \mathcal{H}_{\text{finite}}$, we are assuming that the hypothesis class is realizable. Fix any sample size m . A sample S in this case looks like $\{(x_1, 1), \dots, (x_m, 1)\}$ and an obvious ERM strategy is to output h_A for $A = \{x_1, \dots, x_m\}$. Clearly $L_S(h_A) = 0$ but $L_{\mathcal{D}}(h_A) = 1$.

The previous examples show that the size of the hypothesis class does not characterize whether it is learnable. This characterization is provided by the so-called VC-dimension.

Exercise 6.1

Let \mathcal{H} be a set of functions from \mathcal{X} to $\{0, 1\}$ and let $\mathcal{H}' \subseteq \mathcal{H}$. Assume that $\text{VCdim}(\mathcal{H}') > \text{VCdim}(\mathcal{H})$. Then there exists a set $C \subseteq \mathcal{X}$ that is shattered by \mathcal{H}' but not by \mathcal{H} . This implies that for all $g: C \rightarrow \{0, 1\}$ there exists $h' \in \mathcal{H}'$ such that $g(x) = h'(x)$ for all $x \in C$. Since $h' \in \mathcal{H}$, this implies that \mathcal{H} shatters C , a contradiction.

Exercise 6.2

In this exercise, \mathcal{X} is finite and $k \leq |\mathcal{X}| =: n$.

6.2.1

We claim that

$$\text{VCdim}(\mathcal{H}_{=k}) = \begin{cases} k & \text{if } k \leq \lfloor n/2 \rfloor \\ n - k & \text{if } k > \lfloor n/2 \rfloor \end{cases}$$

Suppose that $k \leq \lfloor n/2 \rfloor$ and consider a subset $C \subset \mathcal{X}$ of size $k + 1$. Then the all-one function on C cannot be extended to a function in $\mathcal{H}_{=k}$ as it maps $k + 1$ elements of \mathcal{X} to 1.

Hence $\text{VCdim}(\mathcal{H}_{=k}) \leq k$. If $|C| = k$ and $g: C \rightarrow \{0, 1\}$ that maps k' elements of C to 1, we can extend g to a function on \mathcal{X} that maps exactly k elements of \mathcal{X} to 1. This shows that $\text{VCdim}(\mathcal{H}_{=k}) \geq k$. Hence $\text{VCdim}\mathcal{H}_{=k} = k$.

Now consider the case $k > \lfloor n/2 \rfloor$. If C is subset of size $n - k + 1$, then the all-zero function on C cannot be extended to a function in $\mathcal{H}_{=k}$. This happens because there are only $n - (n - k + 1) < k$ elements in $\mathcal{X} \setminus C$. Hence $\text{VCdim}(\mathcal{H}_{=k}) \leq n - k$. If $|C| = n - k$ and $g: C \rightarrow \{0, 1\}$ that assigns 1 to k' elements of C , then we can extend g to a function in $\mathcal{H}_{=k}$ as we have at least $k - k'$ elements in $\mathcal{X} \setminus C$ which we can map to 1. This shows that $\text{VCdim}(\mathcal{H}_{=k}) \geq n - k$. Hence $\text{VCdim}(\mathcal{H}_{=k}) = n - k$.

6.2.2

First observe that if $k \geq \lfloor n/2 \rfloor$, then $\mathcal{H}_{\leq k}$ includes all possible functions from \mathcal{X} to $\{0, 1\}$. This is because any function $g: \mathcal{X} \rightarrow \{0, 1\}$ maps at most half the elements of \mathcal{X} to either 0 or 1 and hence is in $\mathcal{H}_{\leq k}$. Hence in this case every subset of \mathcal{X} is shattered by $\mathcal{H}_{\leq k}$ and $\text{VCdim}(\mathcal{H}_{\leq k}) = n$.

If $k < \lfloor n/2 \rfloor$, then we claim that $\text{VCdim}(\mathcal{H}_{\leq k}) = 2k + 1$. Let $C \subset \mathcal{X}$ of size $2k + 1$ and consider a function $g: C \rightarrow \{0, 1\}$. Such a function maps at most k elements to either 0 or 1. Suppose that it maps at most k elements to 1. Extend g to a function on \mathcal{X} by mapping all elements of $\mathcal{X} \setminus C$ to 0. This extension is a function on \mathcal{X} that maps at most k elements to 1 and hence is an element of $\mathcal{H}_{\leq k}$. The reasoning is similar had g mapped at most k elements to 0. This show that $\text{VCdim}(\mathcal{H}_{\leq k}) \geq 2k + 1$.

Now suppose that $C \subset \mathcal{X}$ is of size $2k + 2$. Consider a map that assigns half the elements of C to 0 and the other half to 1. This map cannot be extended to a function in $\mathcal{H}_{\leq k}$. This proves that $\text{VCdim}(\mathcal{H}_{\leq k}) \leq 2k + 1$. Thus:

$$\text{VCdim}(\mathcal{H}_{\leq k}) = \begin{cases} 2k + 1 & \text{if } k < \lfloor n/2 \rfloor \\ n & \text{if } k \geq \lfloor n/2 \rfloor \end{cases}$$

Exercise 6.3

Since $|\mathcal{H}_{n\text{-parity}}| = 2^n$, using the upper bound on the VC-dimension,

$$\text{VCdim}(\mathcal{H}_{n\text{-parity}}) \leq \log_2 |\mathcal{H}_{n\text{-parity}}| = n.$$

We claim that $\text{VCdim}(\mathcal{H}_{n\text{-parity}}) = n$. Let $C = \{c_1, \dots, c_n\} \subset \mathcal{X}$ be the set of standard basis vectors such that c_i is the basis vector with a 1 in the i th position and 0's elsewhere. Let (b_1, \dots, b_n) be a function from C to $\{0, 1\}$. Construct an index set $I \subseteq \{1, \dots, n\}$ as follows: start with $I \leftarrow \emptyset$; for $1 \leq i \leq n$, if $b_i = 1$ then $I \leftarrow I \cup \{i\}$.

We claim that $h_I(c_j) = b_j$ for all $1 \leq j \leq n$. For if $b_j = 0$, then $j \notin I$ and $\sum_{i \in I} c_{ji} = 0 \pmod{2}$; if $b_j = 1$, then $j \in I$ and $\sum_{i \in I} c_{ji} = 1 \pmod{2}$, proving the claim. This shows that every function from C to $\{0, 1\}$ can be extended to a function in $\mathcal{H}_{n\text{-parity}}$. Hence $\text{VCdim}(\mathcal{H}_{n\text{-parity}}) \leq n$ and together with the upper bound for the VC-dimension, this implies that $\text{VCdim}(\mathcal{H}_{n\text{-parity}}) = n$.