

## Learning via Uniform Convergence

The first formal learning model that we have discussed was the PAC model. In Chapter 2 we have shown that under the realizability assumption, any finite hypothesis class is PAC learnable. In this chapter we will develop a general tool, *uniform convergence*, and apply it to show that any finite class is learnable in the agnostic PAC model with general loss functions, as long as the range loss function is bounded.

### 4.1 UNIFORM CONVERGENCE IS SUFFICIENT FOR LEARNABILITY

The idea behind the learning condition discussed in this chapter is very simple. Recall that, given a hypothesis class,  $\mathcal{H}$ , the ERM learning paradigm works as follows: Upon receiving a training sample,  $S$ , the learner evaluates the risk (or error) of each  $h$  in  $\mathcal{H}$  on the given sample and outputs a member of  $\mathcal{H}$  that minimizes this empirical risk. The hope is that an  $h$  that minimizes the empirical risk with respect to  $S$  is a risk minimizer (or has risk close to the minimum) with respect to the true data probability distribution as well. For that, it suffices to ensure that the empirical risks of all members of  $\mathcal{H}$  are good approximations of their true risk. Put another way, we need that uniformly over all hypotheses in the hypothesis class, the empirical risk will be close to the true risk, as formalized in the following.

**Definition 4.1** ( $\epsilon$ -representative sample). A training set  $S$  is called  $\epsilon$ -representative (w.r.t. domain  $Z$ , hypothesis class  $\mathcal{H}$ , loss function  $\ell$ , and distribution  $\mathcal{D}$ ) if

$$\forall h \in \mathcal{H}, \quad |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon.$$

The next simple lemma states that whenever the sample is  $(\epsilon/2)$ -representative, the ERM learning rule is guaranteed to return a good hypothesis.

**Lemma 4.2.** Assume that a training set  $S$  is  $\frac{\epsilon}{2}$ -representative (w.r.t. domain  $Z$ , hypothesis class  $\mathcal{H}$ , loss function  $\ell$ , and distribution  $\mathcal{D}$ ). Then, any output of  $\text{ERM}_{\mathcal{H}}(S)$ , namely, any  $h_S \in \arg\min_{h \in \mathcal{H}} L_S(h)$ , satisfies

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

*Proof.* For every  $h \in \mathcal{H}$ ,

$$L_{\mathcal{D}}(h_S) \leq L_S(h_S) + \frac{\epsilon}{2} \leq L_S(h) + \frac{\epsilon}{2} \leq L_{\mathcal{D}}(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = L_{\mathcal{D}}(h) + \epsilon,$$

where the first and third inequalities are due to the assumption that  $S$  is  $\frac{\epsilon}{2}$ -representative (Definition 4.1) and the second inequality holds since  $h_S$  is an ERM predictor.  $\square$

The preceding lemma implies that to ensure that the ERM rule is an agnostic PAC learner, it suffices to show that with probability of at least  $1 - \delta$  over the random choice of a training set, it will be an  $\epsilon$ -representative training set. The uniform convergence condition formalizes this requirement.

**Definition 4.3** (Uniform Convergence). We say that a hypothesis class  $\mathcal{H}$  has the *uniform convergence property* (w.r.t. a domain  $Z$  and a loss function  $\ell$ ) if there exists a function  $m_{\mathcal{H}}^{\text{UC}} : (0, 1)^2 \rightarrow \mathbb{N}$  such that for every  $\epsilon, \delta \in (0, 1)$  and for every probability distribution  $\mathcal{D}$  over  $Z$ , if  $S$  is a sample of  $m \geq m_{\mathcal{H}}^{\text{UC}}(\epsilon, \delta)$  examples drawn i.i.d. according to  $\mathcal{D}$ , then, with probability of at least  $1 - \delta$ ,  $S$  is  $\epsilon$ -representative.

Similar to the definition of sample complexity for PAC learning, the function  $m_{\mathcal{H}}^{\text{UC}}$  measures the (minimal) sample complexity of obtaining the uniform convergence property, namely, how many examples we need to ensure that with probability of at least  $1 - \delta$  the sample would be  $\epsilon$ -representative.

The term *uniform* here refers to having a fixed sample size that works for all members of  $\mathcal{H}$  and over all possible probability distributions over the domain.

The following corollary follows directly from Lemma 4.2 and the definition of uniform convergence.

**Corollary 4.4.** *If a class  $\mathcal{H}$  has the uniform convergence property with a function  $m_{\mathcal{H}}^{\text{UC}}$  then the class is agnostically PAC learnable with the sample complexity  $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta)$ . Furthermore, in that case, the  $\text{ERM}_{\mathcal{H}}$  paradigm is a successful agnostic PAC learner for  $\mathcal{H}$ .*

## 4.2 FINITE CLASSES ARE AGNOSTIC PAC LEARNABLE

In view of Corollary 4.4, the claim that every finite hypothesis class is agnostic PAC learnable will follow once we establish that uniform convergence holds for a finite hypothesis class.

To show that uniform convergence holds we follow a two step argument, similar to the derivation in Chapter 2. The first step applies the union bound while the second step employs a measure concentration inequality. We now explain these two steps in detail.

Fix some  $\epsilon, \delta$ . We need to find a sample size  $m$  that guarantees that for any  $\mathcal{D}$ , with probability of at least  $1 - \delta$  of the choice of  $S = (z_1, \dots, z_m)$  sampled i.i.d. from  $\mathcal{D}$  we have that for all  $h \in \mathcal{H}$ ,  $|L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$ . That is,

$$\mathcal{D}^m(\{S : \forall h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon\}) \geq 1 - \delta.$$

Equivalently, we need to show that

$$\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) < \delta.$$

Writing

$$\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\} = \cup_{h \in \mathcal{H}} \{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\},$$

and applying the union bound (Lemma 2.2) we obtain

$$\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \leq \sum_{h \in \mathcal{H}} \mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}). \quad (4.1)$$

Our second step will be to argue that each summand of the right-hand side of this inequality is small enough (for a sufficiently large  $m$ ). That is, we will show that for any fixed hypothesis,  $h$ , (which is chosen in advance prior to the sampling of the training set), the gap between the true and empirical risks,  $|L_S(h) - L_{\mathcal{D}}(h)|$ , is likely to be small.

Recall that  $L_{\mathcal{D}}(h) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(h, z)]$  and that  $L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$ . Since each  $z_i$  is sampled i.i.d. from  $\mathcal{D}$ , the expected value of the random variable  $\ell(h, z_i)$  is  $L_{\mathcal{D}}(h)$ . By the linearity of expectation, it follows that  $L_{\mathcal{D}}(h)$  is also the expected value of  $L_S(h)$ . Hence, the quantity  $|L_{\mathcal{D}}(h) - L_S(h)|$  is the deviation of the random variable  $L_S(h)$  from its expectation. We therefore need to show that the measure of  $L_S(h)$  is *concentrated* around its expected value.

A basic statistical fact, the *law of large numbers*, states that when  $m$  goes to infinity, empirical averages converge to their true expectation. This is true for  $L_S(h)$ , since it is the empirical average of  $m$  i.i.d random variables. However, since the law of large numbers is only an asymptotic result, it provides no information about the gap between the empirically estimated error and its true value for any given, finite, sample size.

Instead, we will use a measure concentration inequality due to Hoeffding, which quantifies the gap between empirical averages and their expected value.

**Lemma 4.5** (Hoeffding's Inequality). *Let  $\theta_1, \dots, \theta_m$  be a sequence of i.i.d. random variables and assume that for all  $i$ ,  $\mathbb{E}[\theta_i] = \mu$  and  $\mathbb{P}[a \leq \theta_i \leq b] = 1$ . Then, for any  $\epsilon > 0$*

$$\mathbb{P}\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right] \leq 2 \exp\left(-2m\epsilon^2/(b-a)^2\right).$$

The proof can be found in Appendix B.

Getting back to our problem, let  $\theta_i$  be the random variable  $\ell(h, z_i)$ . Since  $h$  is fixed and  $z_1, \dots, z_m$  are sampled i.i.d., it follows that  $\theta_1, \dots, \theta_m$  are also i.i.d. random variables. Furthermore,  $L_S(h) = \frac{1}{m} \sum_{i=1}^m \theta_i$  and  $L_{\mathcal{D}}(h) = \mu$ . Let us further assume that the range of  $\ell$  is  $[0, 1]$  and therefore  $\theta_i \in [0, 1]$ . We therefore obtain that

$$\mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) = \mathbb{P}\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right] \leq 2 \exp\left(-2m\epsilon^2\right). \quad (4.2)$$

Combining this with Equation (4.1) yields

$$\begin{aligned} \mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) &\leq \sum_{h \in \mathcal{H}} 2 \exp\left(-2m\epsilon^2\right) \\ &= 2|\mathcal{H}| \exp\left(-2m\epsilon^2\right). \end{aligned}$$

Finally, if we choose

$$m \geq \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2}$$

then

$$\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \epsilon\}) \leq \delta.$$

**Corollary 4.6.** *Let  $\mathcal{H}$  be a finite hypothesis class, let  $Z$  be a domain, and let  $\ell : \mathcal{H} \times Z \rightarrow [0, 1]$  be a loss function. Then,  $\mathcal{H}$  enjoys the uniform convergence property with sample complexity*

$$m_{\mathcal{H}}^{UC}(\epsilon, \delta) \leq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2} \right\rceil.$$

Furthermore, the class is agnostically PAC learnable using the ERM algorithm with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\epsilon/2, \delta) \leq \left\lceil \frac{2\log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil.$$

**Remark 4.1** (The “Discretization Trick”). While the preceding corollary only applies to finite hypothesis classes, there is a simple trick that allows us to get a very good estimate of the practical sample complexity of infinite hypothesis classes. Consider a hypothesis class that is parameterized by  $d$  parameters. For example, let  $\mathcal{X} = \mathbb{R}$ ,  $\mathcal{Y} = \{\pm 1\}$ , and the hypothesis class,  $\mathcal{H}$ , be all functions of the form  $h_{\theta}(x) = \text{sign}(x - \theta)$ . That is, each hypothesis is parameterized by one parameter,  $\theta \in \mathbb{R}$ , and the hypothesis outputs 1 for all instances larger than  $\theta$  and outputs  $-1$  for instances smaller than  $\theta$ . This is a hypothesis class of an infinite size. However, if we are going to learn this hypothesis class in practice, using a computer, we will probably maintain real numbers using floating point representation, say, of 64 bits. It follows that in practice, our hypothesis class is parameterized by the set of scalars that can be represented using a 64 bits floating point number. There are at most  $2^{64}$  such numbers; hence the actual size of our hypothesis class is at most  $2^{64}$ . More generally, if our hypothesis class is parameterized by  $d$  numbers, in practice we learn a hypothesis class of size at most  $2^{64d}$ . Applying Corollary 4.6 we obtain that the sample complexity of such classes is bounded by  $\frac{128d + 2\log(2/\delta)}{\epsilon^2}$ . This upper bound on the sample complexity has the deficiency of being dependent on the specific representation of real numbers used by our machine. In Chapter 6 we will introduce a rigorous way to analyze the sample complexity of infinite size hypothesis classes. Nevertheless, the discretization trick can be used to get a rough estimate of the sample complexity in many practical situations.

### 4.3 SUMMARY

If the uniform convergence property holds for a hypothesis class  $\mathcal{H}$  then in most cases the empirical risks of hypotheses in  $\mathcal{H}$  will faithfully represent their true risks. Uniform convergence suffices for agnostic PAC learnability using the ERM rule. We have shown that finite hypothesis classes enjoy the uniform convergence property and are hence agnostic PAC learnable.

#### 4.4 BIBLIOGRAPHIC REMARKS

Classes of functions for which the uniform convergence property holds are also called Glivenko-Cantelli classes, named after Valery Ivanovich Glivenko and Francesco Paolo Cantelli, who proved the first uniform convergence result in the 1930s. See (Dudley, Gine & Zinn 1991). The relation between uniform convergence and learnability was thoroughly studied by Vapnik – see (Vapnik 1992, Vapnik 1995, Vapnik 1998). In fact, as we will see later in Chapter 6, the fundamental theorem of learning theory states that in binary classification problems, uniform convergence is not only a sufficient condition for learnability but is also a necessary condition. This is not the case for more general learning problems (see (Shalev-Shwartz, Shamir, Srebro & Sridharan 2010)).

#### 4.5 EXERCISES

4.1 In this exercise, we show that the  $(\epsilon, \delta)$  requirement on the convergence of errors in our definitions of PAC learning, is, in fact, quite close to a simpler looking requirement about averages (or expectations). Prove that the following two statements are equivalent (for any learning algorithm  $A$ , any probability distribution  $\mathcal{D}$ , and any loss function whose range is  $[0, 1]$ ):

1. For every  $\epsilon, \delta > 0$ , there exists  $m(\epsilon, \delta)$  such that  $\forall m \geq m(\epsilon, \delta)$

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) > \epsilon] < \delta$$

- 2.

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

(where  $\mathbb{E}_{S \sim \mathcal{D}^m}$  denotes the expectation over samples  $S$  of size  $m$ ).

4.2 **Bounded loss functions:** In Corollary 4.6 we assumed that the range of the loss function is  $[0, 1]$ . Prove that if the range of the loss function is  $[a, b]$  then the sample complexity satisfies

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\epsilon/2, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)(b-a)^2}{\epsilon^2} \right\rceil.$$