

# Machine Learning Theory (CSC 482A/581A) - Lectures 6 and 7

Nishant Mehta

## 1 Agnostic Learning

From now on, we will study the agnostic learning setting, wherein the labels themselves can be random conditional on the input. Thus, the distribution  $P$  is now a joint distribution over  $\mathcal{X} \times \mathcal{Y}$ . In PAC learning, it made sense to analyze the rate of convergence of a learning algorithm's risk  $R(\hat{f})$  to zero as the sample size  $n$  increases; however, in the agnostic setting it may no longer be the case that there exists a hypothesis from the set  $\bar{\mathcal{F}}$  of all possible hypotheses that obtains zero risk.<sup>1</sup> A more sensible goal is to hope for a learning algorithm for which the *excess risk* with respect to the best possible hypothesis decays to zero as  $n \rightarrow \infty$ . Let us therefore study the behavior of this best possible hypothesis.

### 1.1 Bayes classifier

**Definition 1.** The Bayes risk  $R^*$  is defined as the minimum risk among all possible hypotheses:<sup>2</sup>

$$R^* = \inf_{f \in \bar{\mathcal{F}}} R(f).$$

A *Bayes optimal classifier*, or Bayes classifier, is a hypothesis  $f$  which obtains the Bayes risk:

$$R(f) = R^*.$$

What form does a Bayes classifier take? For an input  $x \in \mathcal{X}$ , it is easy to see that the conditional risk  $\mathbb{E}[\mathbf{1}[\hat{y} \neq Y] \mid X = x] = \Pr(\hat{y} \neq Y \mid X = x)$  is minimized by predicting

$$\hat{y} \in \arg \max_{y \in \{0,1\}} \Pr(Y = y \mid X = x).$$

Hence, by arbitrarily breaking ties in favor of the positive class, we take the Bayes classifier to be

$$f_{\text{Bayes}}(x) = \mathbf{1}[\Pr(Y = 1 \mid X = x) \geq 1/2].$$

### 1.2 Minimizing excess risk with respect to $f_{\text{Bayes}}$ is hopeless

Now that we have defined the excess risk with respect to  $f_{\text{Bayes}}$ , a natural question arises:

For a fixed input space  $\mathcal{X}$ , is there a learning algorithm  $\mathcal{A}$  for which, no matter the distribution  $P$ , for any  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , there is a sample size  $n(\varepsilon, \delta)$  (not depending on  $P$ ) such that  $R(\hat{f}) - R^* \leq \varepsilon$  with probability at least  $1 - \delta$ ?

<sup>1</sup>Technically, whenever I refer to the set of all possible hypotheses, I actually mean the set of all measurable hypotheses. If you do not know what the term “measurable” means, do not worry about this footnote.

<sup>2</sup>If you do not know what inf means: if you are an undergrad, inf stands for “infimum”, and you may think of it roughly as “minimum”. If you are a graduate student, you should familiarize yourself with infimums and supremums.

Such a learning algorithm would be a universal learner, as it performs well against *any* distribution  $P$ . Unfortunately, this is not possible, as shown by the following “No-Free-Lunch” result:

In agnostic learning, for any learning algorithm and sample size  $n$ , there is a distribution  $P \in \Delta(\mathcal{X} \times \mathcal{Y})$  with deterministic labels under which, with constant positive probability,  $R(\hat{f})$  is lower bounded by a positive constant.

**Remarks:**

- Since the label  $Y$  is deterministic given  $X$ , we have  $P(Y = 1 \mid X) \in \{0, 1\}$ ; consequently, there is a perfect labeling function, and hence  $R(\hat{f}) - R^* = R(\hat{f})$ .

### 1.3 The Agnostic Model

In light of the impossibility result of competing with hypothesis  $f_{\text{Bayes}}$ , we will instead compete against the best hypothesis within our hypothesis space  $\mathcal{F}$ .

Let  $f^*$  be a hypothesis in  $\mathcal{F}$  that minimizes the risk (under distribution  $P$ ), so that

$$R(f^*) = \inf_{f \in \mathcal{F}} R(f).$$

**Definition 2.** We say that  $\mathcal{F}$  is *agnostically learnable* if there exists an algorithm  $\mathcal{A}$  and a function  $n: \mathbb{R}_+ \times (0, 1) \rightarrow \mathbb{N}$  which, for any distribution  $P$  over  $\mathcal{X} \times \mathcal{Y}$  and for all  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , satisfy the following guarantee:

If  $\mathcal{A}$  is given access to  $n(\varepsilon, \delta)$  labeled examples drawn i.i.d. from  $P$ , then with probability at least  $1 - \delta$ ,  $\mathcal{A}$  outputs a hypothesis  $\hat{f}$  with excess risk  $R(\hat{f}) - R(f^*) \leq \varepsilon$ .

We say that  $\mathcal{F}$  is *efficiently agnostically learnable* if, in addition,  $\mathcal{A}$  runs in time polynomial in  $\frac{1}{\varepsilon}$  and  $\frac{1}{\delta}$ .

As with PAC learning, we can also require  $\mathcal{A}$  to output hypotheses in  $\mathcal{F}$  (so that  $\mathcal{A}$  is proper). If  $\mathcal{F}$  is agnostically learnable by such an algorithm, then  $\mathcal{F}$  is *proper agnostically learnable*.

### 1.4 Error decompositions

#### Decomposition of excess risk into approximation error and estimation error

Consider a learning algorithm  $\mathcal{A}$  which, given a training sample  $S$ , outputs some hypothesis  $\hat{f} \in \mathcal{F}$ . The excess risk of  $\hat{f}$  with respect to the Bayes classifier can be decomposed as

$$R(\hat{f}) - R^* = \underbrace{(R(f^*) - R^*)}_{\text{approximation error}} + \underbrace{(R(\hat{f}) - R(f^*))}_{\text{estimation error}}. \quad (1)$$

The first term in the decomposition is the *approximation error*: it is a measure of how well the class  $\mathcal{F}$  can approximate the Bayes classifier in terms of risk; if  $f_{\text{Bayes}} \in \mathcal{F}$ , then the approximation error is zero. Bounding the approximation error requires knowledge of  $R^*$  or some partial information about  $f_{\text{Bayes}}$ , and in settings where little or no distributional assumptions are made, it is thus very difficult if not impossible to control the approximation error. Note, however, that if we begin making certain assumptions about  $P$  and if we allow  $\mathcal{F}$  to grow with  $n$  (so that at sample size  $n$  Learner outputs a hypothesis in  $\mathcal{F}_n$ ), then it is possible to obtain rates of convergence of the approximation error to zero.

The second term in the decomposition is the *estimation error*. Unlike the approximation error, provided that  $\mathcal{F}$  is not “too large” it is possible to obtain good bounds on the estimation error in a distribution-free way, i.e. without having any information about the underlying distribution  $P$ . How does the estimation error typically depend on  $\mathcal{F}$ ? As we will soon see, for learning algorithms that return hypotheses that have low empirical risk, the estimation error increases with  $|\mathcal{F}|$ . This accords with our intuition that, information-theoretically, we need more bits of information to “whittle down”  $\mathcal{F}$  to the risk minimizer (or set of risk minimizers) as  $\mathcal{F}$  increases in size.

The decomposition (1) into approximation error and estimation error highlights the familiar trade-off between *model expressivity* and *generalization*. As we increase the size (complexity) of our model  $\mathcal{F}$ , the approximation error decreases since the model can express more patterns; simultaneously, however, it becomes more likely that we will overfit and hence fail to generalize well.

Our primary focus will be controlling the estimation error. Controlling the estimation error rather than the excess risk with respect to  $f_{\text{Bayes}}$  has various motivations, including

- If we are “lucky” and the approximation error is zero or sufficiently small, a bound on the estimation error also provides a good bound on  $R(\hat{f}) - R^*$ .
- Suppose that we are in a nonparametric setup where, at sample size  $n$ , Learner employs hypothesis space  $\mathcal{F}_n$ . Under mild assumptions about the true distribution, we may be able to control the approximation error as a function of  $n$ . It then is also useful to control the estimation error for each  $\mathcal{F}_n$ , as we then can determine how quickly the complexity of the model should increase with the sample size.

**Oracle inequality approach.** A bound on the estimation error of a learning algorithm  $\mathcal{A}$  that outputs hypothesis  $\hat{f}$  is equivalent to a bound of the form:

$$R(\hat{f}) \leq R(f^*) + \text{BOUND}(\mathcal{F}, n). \quad (2)$$

In statistics and machine learning, a bound of this form is called an *oracle inequality*. The name stems from our comparing the performance of  $\hat{f}$  to that of an omniscient oracle which plays  $f^*$ , the best hypothesis in  $\mathcal{F}$ .

It is natural to seek an oracle inequality for a learning algorithm, as we then know how far off the risk we obtain is from the best possible risk obtainable via  $\mathcal{F}$ . However, to the practitioner, oracle inequalities are not immediately useful:  $R(f^*)$  is an unknown quantity, so, while the bound may be correct, a practitioner has no observable upper bound on  $R(\hat{f})$  (!).

### Deviations approach: Decomposition of risk into empirical risk and deviation

The error decomposition below, this time of the risk of  $R(\hat{f})$  itself, *can* lead to an observable bound. For any hypothesis  $f$  and training sample  $S = ((X_1, Y_1), \dots, (X_n, Y_n))$ , let  $\hat{R}_S(f) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}[f(X_j) \neq Y_j]$  denote the empirical risk of  $f$  on  $S$ . Then

$$R(\hat{f}) = \hat{R}_S(\hat{f}) + \underbrace{(R(\hat{f}) - \hat{R}_S(\hat{f}))}_{\text{deviation}}. \quad (3)$$

Let’s see how we can use (3) to get an upper bound on the risk of  $\hat{f}$ . Suppose that we have a bound of the form:

$$\left| R(f) - \hat{R}_S(f) \right| \leq \varepsilon \quad \text{for all } f \in \mathcal{F}. \quad (4)$$

This bound is known as a uniform deviation bound, since it bounds the deviation of  $\hat{R}_S(f)$  from its mean  $\mathbb{E}[\hat{R}_S(f)] = R(f)$ , uniformly over  $\mathcal{F}$ .

From (4), the bound holds for  $\hat{f}$  in particular, and so we immediately obtain the risk bound

$$R(\hat{f}) \leq \hat{R}_S(\hat{f}) + \varepsilon. \quad (5)$$

Note that this upper bound is observable, since the empirical risk of  $\hat{f}$  can be observed.

Moreover, as we will see below, a bound of the form (4) can, with just a few short steps, lead to an oracle inequality.

## 2 A first excess risk bound for finite classes

Let's derive a first excess risk bound for agnostically learning a finite class  $\mathcal{F}$ . We will obtain a bound by way of a concentration inequality known as Hoeffding's inequality, proved by Wassily Hoeffding in 1963.

**Theorem 1.** *Let  $Z_1, \dots, Z_n$  be independent random variables such that  $Z_j \in [a_j, b_j]$  for  $j \in [n]$ . Let  $\bar{Z} = \frac{1}{n} \sum_{j=1}^n Z_j$ . Then for any  $\varepsilon > 0$ :*

$$\Pr \left( \bar{Z} - \mathbb{E}[\bar{Z}] \geq \varepsilon \right) \leq \exp \left( \frac{-2n^2\varepsilon^2}{\sum_{j=1}^n (b_j - a_j)^2} \right).$$

Before establishing an excess risk bound, we will first establish a uniform convergence result: the empirical risk converges to the actual risk uniformly over  $\mathcal{F}$ . It becomes tiresome to carry around the subscript  $S$  for the empirical risk, so we use the abbreviation  $\hat{R}(f) := \hat{R}_S(f)$ .

**Theorem 2.** *Let  $\mathcal{F}$  be a finite set of hypotheses and let  $P$  be a fixed distribution over  $\mathcal{X} \times \mathcal{Y}$ . For any  $\varepsilon > 0$  and any  $\delta \in (0, 1)$ , if  $(X_1, Y_1), \dots, (X_n, Y_n)$  are drawn i.i.d. from  $P$  with*

$$n \geq \frac{\log |\mathcal{F}| + \log \frac{2}{\delta}}{2\varepsilon^2},$$

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

$$\left| R(f) - \hat{R}(f) \right| \leq \varepsilon \quad \text{for all } f \in \mathcal{F}.$$

*Proof.* Fix some  $f \in \mathcal{F}$  and consider the probability that

$$R(f) - \hat{R}(f) > \varepsilon.$$

This event may be rewritten as

$$\frac{1}{n} \sum_{j=1}^n \mathbf{1}[f(X_j) \neq Y_j] - \mathbb{E} \left[ \frac{1}{n} \sum_{j=1}^n \mathbf{1}[f(X_j) \neq Y_j] \right] > \varepsilon,$$

and so we may apply Hoeffding's inequality twice, once with  $Z_j = -\mathbf{1}[f(X_j) \neq Y_j]$ ,  $a_j = 0$ , and  $b_j = 1$  for  $j \in [n]$ , yielding

$$\Pr \left( R(f) - \hat{R}(f) > \varepsilon \right) \leq e^{-2n\varepsilon^2},$$

and once with once with  $Z_j = \mathbf{1}[f(X_j) \neq Y_j]$ ,  $a_j = -1$ , and  $b_j = 0$  for  $j \in [n]$ , yielding

$$\Pr\left(\hat{R}(f) - R(f) > \varepsilon\right) \leq e^{-2n\varepsilon^2}.$$

Hence,

$$\Pr\left(\left|R(f) - \hat{R}(f)\right| > \varepsilon\right) \leq 2e^{-2n\varepsilon^2}.$$

Next, applying the union bound, we have

$$\begin{aligned} \Pr\left(\exists f \in \mathcal{F} : \left|R(f) - \hat{R}(f)\right| > \varepsilon\right) &\leq \sum_{f \in \mathcal{F}} \Pr\left(\left|R(f) - \hat{R}(f)\right| > \varepsilon\right) \\ &\leq 2|\mathcal{F}|e^{-2n\varepsilon^2}. \end{aligned}$$

The result follows by setting the RHS to  $\delta$  and solving for  $n$ .  $\square$

We now prove that any finite class can be agnostically learned using *empirical risk minimization* (ERM) over  $\mathcal{F}$ , a method which outputs the hypothesis in  $\mathcal{F}$  that minimizes the empirical risk.

**Theorem 3.** *Let  $\mathcal{F}$  be a finite set of hypotheses, let  $P$  be a fixed distribution over  $\mathcal{X} \times \mathcal{Y}$ , and take  $\mathcal{A}$  to be ERM over  $\mathcal{F}$ . For any  $\varepsilon > 0$  and any  $\delta \in (0, 1)$ , if  $\mathcal{A}$  is run on a training sample  $(X_1, Y_1), \dots, (X_n, Y_n)$  drawn i.i.d. from  $P$  with*

$$n \geq \frac{2\left(\log |\mathcal{F}| + \log \frac{2}{\delta}\right)}{\varepsilon^2},$$

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

$$R(\hat{f}) \leq R(f^*) + \varepsilon.$$

*Proof.* First, observe that

$$\begin{aligned} R(\hat{f}) - R(f^*) &= \left(\hat{R}(\hat{f}) + (R(\hat{f}) - \hat{R}(\hat{f}))\right) \\ &\quad - \left(\hat{R}(f^*) + (R(f^*) - \hat{R}(f^*))\right) \\ &= \left(\hat{R}(\hat{f}) - \hat{R}(f^*)\right) + \left(R(\hat{f}) - \hat{R}(\hat{f})\right) + \left(\hat{R}(f^*) - R(f^*)\right) \\ &\leq \left(R(\hat{f}) - \hat{R}(\hat{f})\right) + \left(\hat{R}(f^*) - R(f^*)\right) \\ &\leq 2 \max_{f \in \mathcal{F}} \left|R(f) - \hat{R}(f)\right|, \end{aligned}$$

where the first inequality uses the fact that the empirical risk of ERM is no greater than the empirical risk of  $f^*$ . Next, from Theorem 2, with probability at least  $1 - \delta$

$$\max_{f \in \mathcal{F}} \left|R(f) - \hat{R}(f)\right| \leq \varepsilon/2,$$

and so the result holds.  $\square$