

Formulation of the learning problem, Part 2

Maxim Raginsky

September 11, 2014

1 Agnostic (or model-free) learning

The realizable setting we have focused on in the last lecture rests on certain assumptions, which are not always warranted:

- The assumption that the target concept C^* belongs to \mathcal{C} (or that the target function f^* belongs to \mathcal{F}) means that we are trying to fit a hypothesis to data, which are *a priori* known to have been generated by some member of the model class defined by \mathcal{C} (or by \mathcal{F}). However, in general we may not want to (or be able to) assume much about the data generation process, and instead would like to find the best fit to the data at hand using an element of some model class of our choice.
- The assumption that the training features (or inputs) are labelled noiselessly by $\mathbf{1}_{\{x \in C^*\}}$ (or by $f(x)$) rules out the possibility of noisy measurements or observations.
- Finally, even if the above assumption were true, we would not necessarily have *a priori* knowledge of the concept class \mathcal{C} (or the function class \mathcal{F}) containing the target concept (or function). In that case, the best we could hope for is to pick our own model class and seek the best *approximation* to the unknown target concept (or function) among the elements of that class.

The *model-free learning problem* (also referred to as the *agnostic case*), introduced by Haussler [Hau92], takes a more general decision-theoretic approach and removes the above restrictions. It has the following ingredients:

- Sets X , Y , and U
- A class \mathcal{P} of probability distributions on $Z \triangleq X \times Y$
- A class \mathcal{F} of functions $f : X \rightarrow U$ (the *hypothesis space*)
- A *loss function* $\ell : Y \times U \rightarrow [0, 1]$

The learning process takes place as follows. We obtain an i.i.d. sample $Z^n = (Z_1, \dots, Z_n)$, where each $Z_i = (X_i, Y_i)$ is drawn from the same fixed but unknown $P \in \mathcal{P}$. A *learning algorithm* is a sequence $\mathcal{A} = \{A_n\}_{n=1}^\infty$ of mappings

$$A_n : Z^n \rightarrow \mathcal{F}.$$

As before, let

$$\hat{f}_n = A_n(Z^n) = A_n(Z_1, \dots, Z_n) = A_n((X_1, Y_1), \dots, (X_n, Y_n)).$$

This is the hypothesis emitted by the learning algorithm based on the *training data* Z^n . Note that, by definition, \hat{f}_n is a *random element* of the hypothesis space \mathcal{F} , and that it maps each point $x \in X$ to a point $u = \hat{f}_n(x) \in U$. Following the same steps as in the realizable case, we evaluate the goodness of \hat{f}_n by its expected loss

$$L_P(\hat{f}_n) \triangleq \mathbb{E}_P[\ell(Y, \hat{f}_n(X)) | Z^n] = \int_{X \times Y} \ell(y, \hat{f}_n(x)) P(dx, dy),$$

where the expectation is w.r.t. a random couple $(X, Y) \in Z$ drawn according to the same P but independently of Z^n . Note that $L_P(\hat{f}_n)$ is a random variable since so is \hat{f}_n . In general, we can define the expected risk w.r.t. P for every f in our hypothesis space by

$$L_P(f) \triangleq \mathbb{E}_P[\ell(Y, f(X))] = \int_{X \times Y} \ell(y, f(x)) P(dx, dy)$$

as well as the *minimum risk*

$$L_P^*(\mathcal{F}) \triangleq \inf_{f \in \mathcal{F}} L_P(f).$$

Conceptually, $L_P^*(\mathcal{F})$ is the best *possible* performance of any hypothesis in \mathcal{F} when the samples are drawn from P ; similarly, $L_P(\hat{f}_n)$ is the *actual* performance of the algorithm with access to a training sample of size n . It is clear from definitions that

$$0 \leq L_P^*(\mathcal{F}) \leq L_P(\hat{f}_n) \leq 1.$$

The goal of learning is to guarantee that $L_P(\hat{f}_n)$ is as close as possible to $L_P^*(\mathcal{F})$, whatever the true $P \in \mathcal{P}$ happens to be. In order to speak about this quantitatively, we need to assess the probability of getting a “bad” sample. To that end, we define, similarly to what we have done earlier, the quantity

$$r_{\mathcal{A}}(n, \varepsilon) \triangleq \sup_{P \in \mathcal{P}} P^n(Z^n \in Z^n : L_P(\hat{f}_n) \geq L_P^*(\mathcal{F}) + \varepsilon) \quad (1)$$

for every $\varepsilon > 0$. Thus, a sample $Z^n \sim P^n$ is declared to be “bad” if it leads to a hypothesis whose expected risk on an independent test point $(X, Y) \sim P$ is greater than the smallest possible loss $L_P^*(\mathcal{F})$ by at least ε . We have the following:

Definition 1. We say that a learning algorithm for a problem $(X, Y, U, \mathcal{P}, \mathcal{F}, \ell)$ is PAC to accuracy ε if

$$\lim_{n \rightarrow \infty} r_{\mathcal{A}}(n, \varepsilon) = 0.$$

An algorithm that is PAC to accuracy ε for every $\varepsilon > 0$ is said to be PAC. A learning problem specified by a tuple $(X, Y, U, \mathcal{P}, \mathcal{F}, \ell)$ is model-free (or agnostically) learnable (to accuracy ε) if there exists an algorithm for it which is PAC (to accuracy ε).

Let us look at some examples.

1.1 Function learning in the realizable case

First we show that the model-free framework contains the realizable set-up as a special case. To see this, let X be an arbitrary space and let $Y = U = [0, 1]$. Let \mathcal{F} be a class of functions $f : X \rightarrow [0, 1]$. Let \mathcal{P}_X be a family of probability distributions P_X on X . To each P_X and each $f \in \mathcal{F}$ associate a probability

distribution $P_{X,f}$ on $X \times Y$ as follows: let $X \sim P_X$, and let the conditional distribution of Y given $X = x$ be given by

$$P_{Y|X,f}(B|X = x) = \mathbf{1}_{\{f(x) \in B\}}$$

for all (measurable) sets $B \subseteq Y$. The resulting joint distribution $P_{X,f}$ is then uniquely defined by its action on the “rectangles” $A \times B$, $A \subseteq X$ and $B \subseteq Y$:

$$P_{X,f}(A \times B) \triangleq \int_A P_{Y|X,f}(B|x) P_X(dx) = \int_A \mathbf{1}_{\{f(x) \in B\}} P_X(dx)$$

Finally, let $\mathcal{P} = \{P_{X,f} : f \in \mathcal{F}, P_X \in \mathcal{P}_X\}$. Finally, let $\ell(y, u) \triangleq |y - u|^2$.

Now, fixing a probability distribution $P \in \mathcal{P}$ is equivalent to fixing some $P_X \in \mathcal{P}_X$ and some $f \in \mathcal{F}$. A random element of $Z = X \times Y$ drawn according to such a P has the form $(X, f(X))$, where $X \sim P_X$. An i.i.d. sequence $(X_1, Y_1), \dots, (X_n, Y_n)$ drawn according to P therefore has the form

$$(X_1, f(X_1)), \dots, (X_n, f(X_n)),$$

which is precisely what we had in our discussion of function learning in the realizable case. Next, for any $P = P_{X,f} \in \mathcal{P}$ and any other $g \in \mathcal{F}$, we have

$$\begin{aligned} L_{P_{X,f}}(g) &= \int_{X \times Y} |y - g(x)|^2 P_{X,f}(dx, dy) \\ &= \int_{X \times Y} \mathbf{1}_{\{y=f(x)\}} |y - g(x)|^2 P_X(dx) \\ &= \int_X |f(x) - g(x)|^2 P_X(dx) \\ &= \|f - g\|_{L^2(P_X)}^2, \end{aligned}$$

which is precisely the risk $L_{P_X}(g, f)$ that we have considered in our function learning formulation earlier. Moreover,

$$L_{P_{X,f}}^* = \inf_{g \in \mathcal{F}} L_{P_{X,f}}(g) = \inf_{g \in \mathcal{F}} \|f - g\|_{L^2(P_X)}^2 \equiv 0.$$

Therefore,

$$\begin{aligned} r_{\mathcal{A}}(n, \varepsilon) &= \sup_{P_{X,f} \in \mathcal{P}} P_{X,f}^n \left(Z^n \in Z^n : L_{P_{X,f}}(\hat{f}_n) \geq L_{P_{X,f}}^* + \varepsilon \right) \\ &= \sup_{P_X \in \mathcal{P}_X} \sup_{f \in \mathcal{F}} P_X^n \left(X^n \in X^n : L_P(\hat{f}_n, f) \geq \varepsilon \right) \\ &\equiv \bar{r}_{\mathcal{A}}(n, \varepsilon, \mathcal{P}_X). \end{aligned}$$

Thus, the function learning problem in the realizable case can be covered under the model-free framework as well.

1.2 Learning to classify with noisy labels

Consider the concept learning problem in the realizable case, except that now the labels Y_i , which in the original problem had the form $\mathbf{1}_{\{X_i \in C^*\}}$ for some target concept C^* , are *noisy*. That is, if X_i is a training feature point, then the label Y_i may be “flipped” due to chance, independently of all other X_j 's, $j \neq i$.

The precise formulation of this problem is as follows. Let X be a given feature space, let \mathcal{C} be a concept class on it, and let \mathcal{P}_X be a class of probability distributions on X . Suppose that Nature picks some distribution $P_X \in \mathcal{P}_X$ of the features and some target concept $C^* \in \mathcal{C}$. The training data are generated as follows. First, an i.i.d. sample $X^n = (X_1, \dots, X_n)$ is drawn according to some $P_X \in \mathcal{P}_X$. Then the corresponding labels $Y_1, \dots, Y_n \in \{0, 1\}$ are generated as follows:

$$Y_i = \begin{cases} \mathbf{1}_{\{X_i \in C^*\}}, & \text{with probability } 1 - \eta \\ 1 - \mathbf{1}_{\{X_i \in C^*\}}, & \text{with probability } \eta \end{cases} \quad \text{independently of } X^n, \{Y_j\}_{j \neq i}$$

where $\eta < 1/2$ is the *classification noise rate*.

To cast this problem into the model-free framework, let $Y = U = \{0, 1\}$, let $\mathcal{F} = \{I_C : C \in \mathcal{C}\}$, and let $\ell(y, u) = |y - u|^2$. Define a class \mathcal{P} of probability distributions $\{P_{X,C} : P_X \in \mathcal{P}_X, C \in \mathcal{C}\}$ on $X \times Y = X \times \{0, 1\}$ as follows. Let $X \sim P_X$, and for a given $C \in \mathcal{C}$ consider the conditional probability of $Y = 1$ given $X = x$. If $x \in C$, then $Y = 1$ if and only if there was no error in the label; on the other hand, if $x \notin C$, then $Y = 1$ if and only if there was an error. That is,

$$\begin{aligned} P_{Y|X,C}(1|X = x) &= (1 - \eta)\mathbf{1}_{\{x \in C\}} + \eta\mathbf{1}_{\{x \in C^c\}} \\ &= (1 - \eta)\mathbf{1}_{\{x \in C\}} + \eta(1 - \mathbf{1}_{\{x \in C\}}); \\ P_{Y|X,C}(0|X = x) &= 1 - P_{Y|X,C}(1|X = x) \\ &= \eta\mathbf{1}_{\{x \in C\}} + (1 - \eta)(1 - \mathbf{1}_{\{x \in C\}}). \end{aligned}$$

Then for any measurable set $A \subseteq X$ we will have

$$\begin{aligned} P_{X,C}(A \times \{1\}) &= \int_A P_{Y|X,C}(1|X = x)P_X(\mathrm{d}x) \\ &= \int_A [(1 - \eta)\mathbf{1}_{\{x \in C\}} + \eta(1 - \mathbf{1}_{\{x \in C\}})] P_X(\mathrm{d}x) \\ &= (1 - \eta) \int_A \mathbf{1}_{\{x \in C\}} P_X(\mathrm{d}x) + \eta \int_A P_X(\mathrm{d}x) - \eta \int_A \mathbf{1}_{\{x \in C\}} P_X(\mathrm{d}x) \\ &= \eta P_X(A) + (1 - 2\eta)P_X(A \cap C) \end{aligned} \tag{2}$$

and similarly

$$P_{X,C}(A \times \{0\}) = (1 - \eta)P_X(A) - (1 - 2\eta)P_X(A \cap C). \tag{3}$$

Given a hypothesis $f = I_{C'} \in \mathcal{F}$, we have

$$L_{P_{X,C}}(I_{C'}) = \int_{X \times Y} |y - I_{C'}(x)|^2 P_{X,C}(\mathrm{d}x, \mathrm{d}y).$$

Computing this integral is straightforward but tedious. We start by expanding it as follows:

$$\begin{aligned} &\int_{X \times Y} |y - I_{C'}(x)|^2 P_{X,C}(\mathrm{d}x, \mathrm{d}y) \\ &= \int_X |0 - I_{C'}(x)|^2 P_{X,C}(\mathrm{d}x \times \{0\}) + \int_X |1 - I_{C'}(x)|^2 P_{X,C}(\mathrm{d}x \times \{1\}) \\ &= \int_X \mathbf{1}_{\{x \in C'\}} P_{X,C}(\mathrm{d}x \times \{0\}) + \int_X \mathbf{1}_{\{x \in (C')^c\}} P_{X,C}(\mathrm{d}x \times \{1\}) \\ &= P_{X,C}(C' \times \{0\}) + P_{X,C}((C')^c \times \{1\}). \end{aligned}$$

Substituting the expressions (2) and (3) into the above, we get

$$\begin{aligned}
L_{P_{X,C}}(I_{C'}) &= (1-\eta)P_X(C') - (1-2\eta)P_X(C \cap C') + \eta P_X((C')^c) + (1-2\eta)P_X(C \cap (C')^c) \\
&= (1-\eta)\underbrace{(P_X(C \cap C') + P_X(C^c \cap C'))}_{P_X(C')} - (1-2\eta)P_X(C \cap C') \\
&\quad + \eta\underbrace{(P_X(C \cap (C')^c) + P_X(C^c \cap (C')^c))}_{P_X((C')^c)} + (1-2\eta)P_X(C \cap (C')^c) \\
&= (1-\eta)\underbrace{(P_X(C \cap (C')^c) + P_X(C^c \cap C'))}_{P_X(C \Delta C')} + \eta P_X(C \cap C') + \eta\underbrace{P_X(C^c \cap (C')^c)}_{P_X((C \cup C')^c)} \\
&= (1-\eta)P_X(C \Delta C') + \eta P_X(C \cap C') + \eta(1 - P_X(C \cup C')) \\
&= (1-\eta)P_X(C \Delta C') + \eta - \eta\underbrace{(P_X(C \cup C') - P_X(C \cap C'))}_{P_X(C \Delta C')} \\
&= \eta + (1-2\eta)P_X(C \Delta C') \\
&\equiv \eta + (1-2\eta)L_{P_X}(C', C).
\end{aligned}$$

From this, we have

$$\begin{aligned}
L_{P_{X,C}}^*(\mathcal{F}) &= \inf_{C' \in \mathcal{C}} L_{P_{X,C}}(I_{C'}) \\
&= \eta + (1-2\eta) \inf_{C' \in \mathcal{C}} P_X(C \Delta C') \\
&= \eta,
\end{aligned}$$

where the infimum is achieved by letting $C' = C$. From this it follows that

$$L_{P_{X,C}}(C') \geq L_{P_{X,C}}^* + \varepsilon \iff P_{X,C}(C \Delta C') \geq \frac{\varepsilon}{1-2\eta}$$

In other words, learning a concept to accuracy ε with noise rate η is equivalent to learning a concept to accuracy $\varepsilon/(1-2\eta)$ in the noise-free case:

$$r_{\mathcal{A}}(n, \varepsilon) = \bar{r}\left(n, \frac{\varepsilon}{1-2\eta}, \mathcal{P}_X\right).$$

2 Empirical risk minimization

Having formulated the model-free learning problem, we must now turn to the question of how to construct PAC learning algorithms (and the related question of when a hypothesis class is PAC-learnable in the model-free setting).

We will first start with a heuristic argument and then make it rigorous. Suppose we are faced with the learning problem specified by $(X, Y, U, \mathcal{P}, \mathcal{F}, \ell)$. Given a training set $Z^n = (Z_1, \dots, Z_n)$, where each $Z_i = (X_i, Y_i)$ is independently drawn according to some unknown $P \in \mathcal{P}$, what should we do? The first thing to note is that, for any hypothesis $f \in \mathcal{F}$, we can approximate its risk $L_P(f)$ by the *empirical risk*

$$\frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)), \quad (4)$$

whose expectation w.r.t. the distribution of Z^n is clearly equal to $L_P(f)$. In fact, since ℓ is bounded between 0 and 1, Hoeffding's inequality tells us that

$$\left| \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)) - L_P(f) \right| < \varepsilon \quad \text{with probability at least } 1 - 2e^{-2n\varepsilon^2}.$$

We can express these statements more succinctly if we define, for each $f \in \mathcal{F}$, the function $\ell_f : Z \rightarrow [0, 1]$ by

$$\ell_f(z) \equiv \ell_f(x, y) \triangleq \ell(y, f(x)). \quad (5)$$

Then the empirical risk (4) is just the expectation of ℓ_f w.r.t. the empirical distribution P_{Z^n} :

$$P_{Z^n}(\ell_f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)),$$

and, since $L_P(f) = \mathbb{E}_P[\ell(Y, f(X))] = P(\ell_f)$, we will have

$$|P_{Z^n}(\ell_f) - P(\ell_f)| < \varepsilon \quad \text{with probability at least } 1 - 2e^{-2n\varepsilon^2}. \quad (6)$$

Now, given the data Z^n we can compute the empirical risks $P_{Z^n}(\ell_f)$ for every f in our hypothesis class \mathcal{F} . Since (6) holds for each $f \in \mathcal{F}$ individually, we may intuitively claim that the empirical risk for each f is a sufficiently accurate estimator of the corresponding true risk $L_P(f) \equiv P(\ell_f)$. Thus, a reasonable learning strategy would be to find any $\hat{f}_n \in \mathcal{F}$ that would *minimize* the empirical risk, i.e., take

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{F}} P_{Z^n}(\ell_f) = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)). \quad (7)$$

The reason why we would expect something like (7) to work is as follows: if a given f^* is a minimizer of $L_P(f) = P(\ell_f)$ over \mathcal{F} ,

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} P(\ell_f),$$

then its empirical risk, $P_{Z^n}(f^*)$, will be close to $L_P(f^*) = P(\ell_{f^*}) = L_P^*(\mathcal{F})$ with high probability. Moreover, it makes sense to expect that, in some sense, \hat{f}_n defined in (7) would be “close” to f^* , resulting in something like

$$P(\hat{f}_n) \approx P_{Z^n}(\hat{f}_n) \approx P_{Z^n}(f^*) \approx P(f^*)$$

with high probability.

Unfortunately, this is not true in general. However, as we will now see, it is true under certain regularity conditions on the objects \mathcal{P} , \mathcal{F} , and ℓ . In order to state these regularity conditions precisely, let us define the *induced loss function class*

$$\mathcal{L}_{\mathcal{F}} \triangleq \{\ell_f : f \in \mathcal{F}\}.$$

Each $\ell_f \in \mathcal{L}_{\mathcal{F}}$ corresponds to the hypothesis $f \in \mathcal{F}$ via (5). Now, for any $n \in \mathbb{N}$ and any $\varepsilon > 0$ let us define

$$q(n, \varepsilon) \triangleq \sup_{P \in \mathcal{P}} P^n \left(Z^n \in Z^n : \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon \right). \quad (8)$$

For a fixed $P \in \mathcal{P}$, quantity $\sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)|$ is the *worst-case deviation* between the empirical means $P_{Z^n}(\ell_f)$ and their expectations $P(\ell_f)$ over the entire hypothesis class \mathcal{F} . Given P , we say that an i.i.d. sample $Z^n \in \mathcal{Z}^n$ is “bad” if there exists at least one $f \in \mathcal{F}$, for which

$$|P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon.$$

Equivalently, a sample is bad if

$$\sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon.$$

The quantity $q(n, \varepsilon)$ then compensates for the fact that P is unknown by considering the *worst case* over the entire class \mathcal{P} . With this in mind, we make the following definition:

Definition 2. We say that the induced class $\mathcal{L}_{\mathcal{F}}$ has the uniform convergence of empirical means (UCEM) property w.r.t. \mathcal{P} if

$$\lim_{n \rightarrow \infty} q(n, \varepsilon) = 0$$

for every $\varepsilon > 0$.

Theorem 1. If the induced class $\mathcal{L}_{\mathcal{F}}$ has the UCEM property, then the empirical risk minimization (ERM) algorithm of (7) is PAC.

Proof. Fix $\varepsilon, \delta > 0$. We will now show that we can find a sufficiently large $n(\varepsilon, \delta)$, such that $r_{\mathcal{A}}(n, \varepsilon) < \delta$ for all $n \geq n(\varepsilon, \delta)$, where $r_{\mathcal{A}}(n, \varepsilon)$ is defined in (1).

Let $f^* \in \mathcal{F}$ minimize the true risk w.r.t. P , i.e., $P(f^*) = L_P^*(\mathcal{F})$. For any n , we have

$$\begin{aligned} L_P(\hat{f}_n) - L_P^* &= P(\ell_{\hat{f}_n}) - P(f^*) \\ &= \underbrace{P(\ell_{\hat{f}_n}) - P_{Z^n}(\ell_{\hat{f}_n})}_{T_1} + \underbrace{P_{Z^n}(\ell_{\hat{f}_n}) - P_{Z^n}(\ell_{f^*})}_{T_2} + \underbrace{P_{Z^n}(\ell_{f^*}) - P(\ell_{f^*})}_{T_3}, \end{aligned}$$

where in the second line we have added and subtracted $P_{Z^n}(\ell_{\hat{f}_n})$ and $P_{Z^n}(\ell_{f^*})$. We will now analyze the behavior of the three terms, T_1 , T_2 , and T_3 . Since \hat{f}_n minimizes the empirical risk $P_{Z^n}(\ell_f)$ over all $f \in \mathcal{F}$, we will have

$$T_2 = P_{Z^n}(\ell_{\hat{f}_n}) - P_{Z^n}(\ell_{f^*}) \leq 0.$$

Next,

$$T_1 = P(\ell_{\hat{f}_n}) - P_{Z^n}(\ell_{\hat{f}_n}) \leq \sup_{f \in \mathcal{F}} [P_{Z^n}(\ell_f) - P(\ell_f)] \leq \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)|,$$

and the same upper bound holds for T_3 . Hence,

$$L_P(\hat{f}_n) - L_P^*(\mathcal{F}) \leq 2 \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)|. \quad (9)$$

Now, since $\mathcal{L}_{\mathcal{F}}$ has the UCEM property, we can find some sufficiently large $n_0(\varepsilon, \delta)$, such that

$$q(n, \varepsilon/2) = \sup_{P \in \mathcal{P}} P^n \left(Z^n \in \mathcal{Z}^n : \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon/2 \right) < \delta, \quad \forall n \geq n_0(\varepsilon, \delta).$$

From this it follows that, for all $n \geq n_0(\varepsilon, \delta)$, we will have

$$P^n \left(Z^n : \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon/2 \right) < \delta, \quad \forall P \in \mathcal{P}.$$

From (9), we see that

$$L_P(\widehat{f}_n) \geq L_P^*(\mathcal{F}) + \varepsilon \quad \implies \quad \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon/2$$

for all n . However, for all $n \geq n_0(\varepsilon, \delta)$ the latter event will occur with probability at most δ , no matter which P is in effect. Therefore, for all $n \geq n_0(\varepsilon, \delta)$ we will have

$$\begin{aligned} r_{\mathcal{A}}(n, \varepsilon) &= \sup_{P \in \mathcal{P}} P^n \left(Z^n : L_P(\widehat{f}_n) \geq L_P^*(\mathcal{F}) + \varepsilon \right) \\ &\leq \sup_{P \in \mathcal{P}} P^n \left(Z^n : \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \geq \varepsilon/2 \right) \\ &\equiv q(n, \varepsilon/2) \\ &< \delta, \end{aligned}$$

which is precisely what we wanted to show. Thus, $r_{\mathcal{A}}(n, \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$ for every $\varepsilon > 0$, which means that the ERM algorithm is PAC. \square

This theorem shows that the UCEM property of the induced class $\mathcal{L}_{\mathcal{F}}$ is a sufficient condition for the ERM algorithm to be PAC. Now the whole affair rests on us being able to establish the UCEM property for various “interesting” and “useful” problem specifications. This will be our concern in the lectures ahead. However, let me give you a hint of what to expect. In many cases, we will be able to show that the induced class $\mathcal{L}_{\mathcal{F}}$ is so well-behaved that the bound

$$\mathbb{E}_{P^n} \left[\sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)| \right] \leq \frac{C_{\mathcal{F}, \ell}}{\sqrt{n}} \quad (10)$$

holds for every P , where $C_{\mathcal{F}, \ell} > 0$ is some constant that depends only on the characteristics of the hypothesis class \mathcal{F} and the loss function ℓ . Since ℓ_f is bounded between 0 and 1, the function

$$g(Z^n) \triangleq \sup_{f \in \mathcal{F}} |P_{Z^n}(\ell_f) - P(\ell_f)|$$

has bounded differences with constants $c_1 = \dots = c_n = 1/n$. McDiarmid’s inequality then tells us that, for any $t > 0$,

$$P^n \left(g(Z^n) - \mathbb{E}g(Z^n) \geq t \right) \leq e^{-2nt^2}. \quad (11)$$

Let

$$n_0(\varepsilon, \delta) \triangleq \max \left\{ \frac{4C_{\mathcal{F}, \ell}^2}{\varepsilon^2}, \frac{2}{\varepsilon^2} \log \left(\frac{1}{\delta} \right) \right\} + 1. \quad (12)$$

Then for any $n \geq n_0(\varepsilon, \delta)$

$$\begin{aligned} P^n \left(g(Z^n) \geq \varepsilon \right) &= P^n \left(g(Z^n) - \mathbb{E}g(Z^n) \geq \varepsilon - \mathbb{E}g(Z^n) \right) \\ &\leq P^n \left(g(Z^n) - \mathbb{E}g(Z^n) \geq \varepsilon - \frac{C_{\mathcal{F}, \ell}}{\sqrt{n}} \right) && \text{because of (10)} \\ &\leq P^n \left(g(Z^n) - \mathbb{E}g(Z^n) \geq \frac{\varepsilon}{2} \right) && \text{because } n > \frac{4C_{\mathcal{F}, \ell}^2}{\varepsilon^2} \\ &\leq e^{-n\varepsilon^2/2} && \text{using (11) with } t = \varepsilon/2 \\ &< \delta && \text{because } n > \frac{2}{\varepsilon^2} \log \left(\frac{1}{\delta} \right) \end{aligned}$$

for *any* probability distribution P over $Z = X \times Y$. Thus, we have derived a very important fact: If the induced loss class $\mathcal{L}_{\mathcal{F}}$ satisfies (10), then (a) it has the UCEM property, and consequently is model-free learnable using the ERM algorithm, and (b) the sample complexity is polynomial in $1/\epsilon$ and logarithmic in $1/\delta$. Our next order of business will be to derive sufficient conditions on \mathcal{F} and ℓ for something like (10) to hold.

References

- [Hau92] D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 95:129–161, 1992.