

Lesson 5

Probably Approximately Correct

5.1 Probably Approximately Correct

What follows below is taken verbatim from [Kearns and Vazirani \(1994, §1.2\)](#). I mostly just omitted sentences which referred to previous examples.

5.1.1 Definition of the PAC Model

Let X be a set called the **instance space**. We think of X as being a set of encodings of instances or objects in the learner's world. In our rectangle game, the instance space X was simply the set of all points in the Euclidean plane \mathfrak{R}^2 . As another example, in a character recognition application, the instance space might consist of all 2-dimensional arrays of binary pixels of a given width and height.

A **concept** over X is just a subset $c \subseteq X$ of the instance space. In the rectangle game, the concepts were axis-aligned rectangular regions.

A concept c can thus be thought of as the set of all instances that positively exemplify some simple or interesting rule. We can equivalently define a concept to be a boolean mapping $c : X \rightarrow \{0, 1\}$, with $c(x) = 1$ indicating that x is a positive example of c and $c(x) = 0$ indicating that x is a negative example. For this reason, we also sometimes call X the **input space**.

A **concept class** \mathcal{C} over X is a collection of concepts over X .

In our model, a learning algorithm will have access to positive and negative examples of an unknown **target concept** c , chosen from a known concept class \mathcal{C} . The learning algorithm will be judged by its ability to identify a hypothesis concept that can accurately classify instances as positive or negative examples of c . Before specifying the learning protocol further, it is important to note that in our model, learning algorithms "know" the target class \mathcal{C} , in the sense that the designer of the learning algorithm is guaranteed that the target concept will be chosen from \mathcal{C} (but must design the algorithm to work for any $c \in \mathcal{C}$). Let \mathcal{D} be any fixed probability distribution over the instance space X .

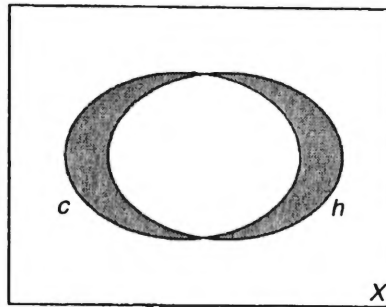


Figure 5.1: Figure 1.4 Venn diagram of two concepts, with symmetric difference shaded.

We will refer to \mathcal{D} as the **target distribution**. If h is any concept over X , then the distribution \mathcal{D} provides a natural measure of **error** between h and the target concept c : namely, we define

$$\text{error}(h) = \Pr_{x \in \mathcal{D}} [c(x) \neq h(x)].$$

Here we regard the concepts c and h as boolean functions, and we have introduced a notational convention that we shall use frequently: the subscript $x \in \mathcal{D}$ to $\Pr[\cdot]$ indicates that the

probability is taken with respect to the random draw of x according to \mathcal{D} . Note that $error(h)$ has an implicit dependence on c and \mathcal{D} that we will usually omit for brevity when no confusion will result.

A useful alternative way to view $error(h)$ is represented in Figure 1.4. Here we view the concepts c and h as sets rather than as functions, and we have drawn an abstract Venn diagram showing the positive examples of c and h , which of course lie within the entire instance space X . Then $error(h)$ is simply the probability with respect to \mathcal{D} that an instance is drawn falling in the shaded region.

Let $EX(c, \mathcal{D})$ be a procedure (we will sometimes call it an *oracle*) that runs in unit time, and on each call returns a labeled example $(x, c(x))$ where x is drawn randomly and independently according to \mathcal{D} . A learning algorithm will have access to this oracle when learning the target concept $c \in C$. Ideally, the learning algorithm will satisfy three properties:

1. The number of calls to $EX(c, \mathcal{D})$ is small, in the sense that it is bounded by a fixed polynomial in some parameters to be specified shortly.
2. The amount of computation performed is small.
3. The algorithm outputs a hypothesis concept h such that $error(h)$ is small.

Note that the number of calls made by a learning algorithm to $EX(c, \mathcal{D})$ is bounded by the running time of the learning algorithm.

We are now ready to give the definition of Probably Approximately Correct learning. We designate it as our preliminary definition, since we shall soon make some important additions to it.

Definition 1 (The PAC Model, Preliminary Definition) Let C be a concept class over X . We say that C is PAC learnable if there exists an algorithm L with the following property: for every concept $c \in C$, for every distribution V on X , and for all $0 < \epsilon < 1/2$ and $0 < \delta < 1/2$, if L is given access to $EX(c, V)$, L outputs a hypothesis concept $h \in C$ satisfying $error(h) < \epsilon$. This probability is taken over the random examples drawn by calls to $EX(c, V)$, and any internal randomization of L .

If L runs in time polynomial in $1/\epsilon$ and $1/\delta$, we say that C is efficiently PAC learnable. We will sometimes refer to the input ϵ as the error parameter, and the input δ as the confidence parameter.

5.2 Working PAC definition

Consider a concept class $C \subseteq P(X)$. C is PAC learnable iff there exists a learning algorithm L with the following property:

- for all $c \in C$,
- for all D over X ,
- for all $0 < \epsilon, \delta$,
- there exists $m = f(\frac{1}{\epsilon}, \frac{1}{\delta})$ such that
- f is a polynomial function, and
- after drawing m samples from $EX(c, D)$, the probability that L outputs a hypothesis h with $\text{error}_{c,D}(h) < \epsilon$ is at least $1 - \delta$.

(Later this definition is modified so that $m = f(\text{size}(h), \frac{1}{\epsilon}, \frac{1}{\delta})$ where $\text{size}(h)$ is a measure of the size of the representation of the concept h .)

5.3 PAC analysis of Axis-Aligned Rectangles

Here we prove that the axis-aligned rectangles is PAC learnable by the rectangle learning strategy discussed by [Kearns and Vazirani \(1994, chapter 1\)](#).

Recall that $\text{error}(h) = \Pr_{x \in D}[c(x) \neq h(x)]$. We want to bound this error by ϵ with probability δ . How do we keep the error smaller than ϵ ? The only area that contributes to the error is the region between the target R and the hypothesized rectangle R' , which is $R \setminus R'$. So we want the probability associated with $R \setminus R'$ to be less than ϵ .

We divide $R \setminus R'$ into four overlapping strips and consider one strip T' . We want to bound the error associated with T' to be less than $\epsilon/4$ so we can ultimately be sure the whole area will have error less than ϵ .

Consider the strip T whose error under the distribution D equals $\epsilon/4$. If T' properly includes T then the error associated with T' exceeds $\epsilon/4$. If this happens, $\Pr[\text{error}(h)]$ could be greater than ϵ . So we want to show that the error of T' is bounded by the error associated with T , which equals $\epsilon/4$.

Note that if any point in T appears in S then in fact T includes T' . This is because if a point in T occurs in S then T' only extends as deep as that point since R' includes all positive points. And if T includes T' then the error associated with T' is bounded by $\epsilon/4$.

What is the probability that a point in S is in T (from which it would ultimately follow that $\Pr[\text{error}(h) < \epsilon]$)? The region T is defined to be the probability that a random draw from $EX(c, D)$ lies in T is $\epsilon/4$. Therefore the probability that a random draw from $EX(c, D)$ does *not* lie in T is $1 - \epsilon/4$. It follows that the probability that none of m random draws lie in T is $(1 - \epsilon/4)^m$. (From which it will ultimately follow that $\Pr[\text{error}(h) > \epsilon]$ is bounded by that number.)

Since there are 4 overlapping strips like T , the probability that none of m random draws lies in any of the 4 strips is less than $4(1 - \epsilon/4)^m$. Formally, we have established

$$\Pr[\text{error}(h) > \epsilon] < 4(1 - \epsilon/4)^m.$$

We want this probability to be less than δ :

$$\Pr[\text{error}(h) > \epsilon] < 4(1 - \epsilon/4)^m < \delta.$$

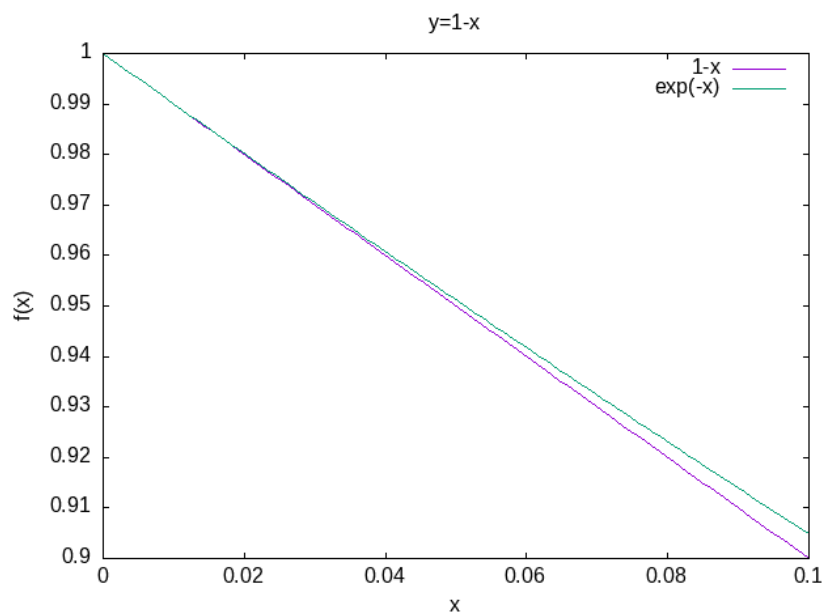
In the region $[0,1]$, it is a fact that

$$1 - x \leq e^{-x} .$$

It follows that

$$(1 - x)^m \leq e^{-mx}$$

in the same region.



Consequently we have shown the following.

$$\Pr[\text{error}(h) > \epsilon] < 4(1 - \epsilon/4)^m \leq 4e^{-m\epsilon/4}$$

Thus any value of m which satisfies

$$4e^{-m\epsilon/4} < \delta$$

will also satisfy

$$\Pr[\text{error}(h) > \epsilon] < \delta ,$$

which is equivalent to

$$\Pr[\text{error}(h) < \epsilon] > 1 - \delta .$$

Here is a complete derivation for finding such m .

$$\begin{aligned}
 4e^{-m\epsilon/4} &< \delta \\
 e^{-m\epsilon/4} &< \delta/4 \\
 -m\epsilon/4 &< \ln(\delta/4) \\
 -m\epsilon/4 &< \ln(\delta) - \ln(4) \\
 m\epsilon/4 &> \ln(4) - \ln(\delta) \\
 m\epsilon/4 &> \ln(4/\delta) \\
 m &> (4/\epsilon) \ln(4/\delta)
 \end{aligned}$$

[Kearns and Vazirani \(1994, p. 6\)](#) write “In summary, provided our tightest-fit algorithm takes a sample of at least $(4/\epsilon) \ln(4/\delta)$ examples to form its hypothesis rectangle R' , we can assert that with probability at least $1 - \delta$, R' will misclassify a new point (drawn according to the same distribution from which the sample was chosen) with probability at most ϵ .”

5.4 Monomials

Monomials is another term for the conjunctions of literals. The elimination algorithm [Valiant \(2013\)](#) discusses in Chapter 5 PAC-learns monomials and in this section we see mathematically why that is the case. But first let’s review what the concept class is, and how the elimination algorithm proceeds.

5.4.1 Variables and Literals

A literal is either positive (x) or negative (\bar{x} , also sometimes written $\neg x$) variable. If there are n variables then there are $2n$ literals. The variables can refer to any property such “has eyes,” “more than 100 pounds,” or “contains a stressed syllable”. The instance space X contains elements a which can be evaluated according to these variables. If a has property x then the positive literal x is true of a otherwise it is false. Conversely, if a does not have the property x then the positive literal x is false of a otherwise it is true.

Suppose we have n variables: x_1, x_2, \dots, x_n . The conjunction of literals then is a formula like the following:

$$x_1 \wedge \bar{x}_2 \wedge x_4 .$$

This means “Elements in X satisfying the formula have property x_1 , do not have property x_2 and have property x_4 .” So each possible formula defines a concept, and every finite set of variables defines a concept class by considering all possible conjunctions of positive and negative literals.

An example familiar to students from computational linguistics 2 would be the Strictly 2-Local languages. These are formal languages that can be described by forbidding finitely many substrings of size 2. For example, if $\Sigma = \{a, b\}$ then the “baba” language $\{ba, baba, bababa, \dots\}$ is given by forbidding the substrings $\times a, bb, aa, b\times, \times\times$.

The variables are all the substrings of length 2 drawn from $\{\bowtie\}\Sigma^*\{\bowtie\}$. Each positive literal x is interpreted as “contains the substring x ” and negative literal \bar{x} is interpreted as “does not contain the substring x .” Thus a monomial describing the “baba” language is shown below.

$$\overline{\bowtie a} \wedge \overline{bb} \wedge \overline{aa} \wedge \overline{b\bowtie} \wedge \overline{\bowtie\bowtie}.$$

5.4.2 The Elimination Algorithm

1. Set $h = x_1 \wedge \bar{x}_1 \wedge x_2 \wedge \bar{x}_2 \wedge \dots \wedge x_n \wedge \bar{x}_n$
2. Receive $(a, c(a))$ from $EX(c, D)$. If a is a negative example (so $c(a) = 0$) repeat this step; otherwise move on.
3. For each $1 \leq i \leq n$: if x_i is true of a then remove \bar{x}_i from h and if \bar{x}_i is true of a then remove x_i from h . Return to step 2.

Note this process never ends!

Here is an example using the “baba” language. Below is a list of all possible literals of length 2 from $\{\bowtie\}\Sigma^*\{\bowtie\}$.

$$\begin{array}{cccccccccc} \bowtie a & \bowtie b & a\bowtie & b\bowtie & aa & ab & ba & bb & \bowtie\bowtie \\ \hline \overline{\bowtie a} & \overline{\bowtie b} & \overline{a\bowtie} & \overline{b\bowtie} & \overline{aa} & \overline{ab} & \overline{ba} & \overline{bb} & \overline{\bowtie\bowtie} \end{array}$$

If the first positive example returned by $EX(c, D)$ is “ba” (so with word boundaries $\bowtie ba \bowtie$), then the following literals would be removed from h :

$$\bowtie a, \overline{\bowtie b}, \overline{a\bowtie}, b\bowtie, aa, ab, \overline{ba}, bb, \bowtie\bowtie$$

5.4.3 PAC analysis

Like the axis-aligned rectangles, the elimination algorithm only ever considers hypotheses h that cover all the observed positive examples. Furthermore, the literals in h always include the literals in the target c because all the literals are in h at the beginning and they are only removed when they contradict c . Because h includes the literals in c , any negative example will always be classified as negative by h . In other words, h never errs on negative examples.

So a literal (positive or negative) z in h only causes h to err on positive examples a where z is not true of a . The total probability of z not being true of a positive example a is denoted by $p(z)$ defined below.

$$p(z) = \Pr_{a \in D} [c(a) = 1, z \text{ is not true of } a]$$

It follows that

$$\text{error}(h) \leq \sum_{z \in h} p(z).$$

We would like to bound $\text{error}(h)$ by ϵ . Since there are at most $2n$ literals z in h , we can achieve this if, for all z , $p(z) \leq \epsilon/2n$.

Call z a bad literal if $p(z) \geq \epsilon/2n$. Consider some bad literal z . The probability that z is removed from h after m calls to $EX(c, D)$ is $(1 - p(z))^m$. Since $p(z)$ is at least $\epsilon/2n$ then this probability is at most $(1 - \epsilon/2n)^m$. Since there are at most $2n$ bad literals, the probability that some bad literal is not removed from h is at most $2n(1 - \epsilon/2n)^m$.

In other words, $\Pr[\text{error}(h) > \epsilon]$ is bounded by $2n(1 - \epsilon/2n)^m$. So if we can find values of m that bound this latter value by δ , then we can conclude that $\Pr[\text{error}(h) < \epsilon] > 1 - \delta$.

Here is a complete derivation for finding such m .

$$\begin{aligned}
 2ne^{-m\epsilon/2n} &< \delta \\
 e^{-m\epsilon/2n} &< \delta/2n \\
 -m\epsilon/2n &< \ln(\delta/2n) \\
 -m\epsilon/2n &< \ln(\delta) - \ln(2n) \\
 m\epsilon/2n &> \ln(2n) - \ln(\delta) \\
 m\epsilon/2n &> \ln(2n/\delta) \\
 m &> (2n/\epsilon) \ln(2n/\delta)
 \end{aligned}$$

Since m is a polynomial function in terms of n , δ , and ϵ , we see that the elimination algorithm PAC-learns the concepts expressible as the conjunction of positive and negative literals.

In the case of our example, $n = 9$. So if we want to be 99% confident that the probability of the error is less than 1% then

$$m = 18/.01 \ln(18/.01) = 1800 \times 7.496 \approx 13492$$

examples suffice for any concept $c \in C$, and for any distribution D over X .

Exercise 1. What if $\Sigma = \{a, b, c\}$?

Bibliography

Kearns, Michael, and Umesh Vazirani. 1994. *An Introduction to Computational Learning Theory*. MIT Press.

Valiant, Leslie. 2013. *Probably Approximately Correct: Nature's Algorithms for Learning and Prospering in a Complex World*. Basic Books.