

Probably Almost Correct (PAC) Learning

Anca Ralescu
MLCI Lab, EECS Dept.
University of Cincinnati

PAC learning

Notation

1. \mathbf{X} the set of examples (training and test). An example is a vector $\mathbf{x} \in \mathbf{X}$.
2. $\{-1, +1\}$ is the set of labels (binary classification);
3. $\mathbf{X} \times \{-1, +1\}$ the set of input/output pairings;
4. P (unknown) probability distribution on \mathbf{X} such that $\mathbf{x} \in \mathbf{X}$ is distributed according to P (we say that the examples are iid - identically independent distributed according to P which means that “the probability of selecting an example (\mathbf{x}, y) , with $y \in \{-1, +1\}$ is the same for any (\mathbf{x}, y) and is independent on the selection of any other such pair”).

Let $f_\alpha(\mathbf{x})$ denote a family of functions (that is, a learning machine) and let $err_P(f_\alpha)$ denote the following quantity:

$$\begin{aligned} err_P(f_\alpha) &= P(\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}) \\ &= \int I_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}} dP(\mathbf{x}, y) \\ &= E_P(I_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}}) \end{aligned} \tag{1}$$

In equation (1) I_A denotes the indicator function (event) defined as $I_A(a) = 1$ if $a \in A$, 0 otherwise.

We show now that $err_P(f_\alpha)$ is the same as $R(\alpha)$.

Indeed,

$$\begin{aligned} R(\alpha) &= \int \frac{1}{2} |f_\alpha(\mathbf{x}) - y| dP(\mathbf{x}, y) \\ &= \int_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}} \frac{1}{2} |f_\alpha(\mathbf{x}) - y| dP(\mathbf{x}, y) + \int_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) = y\}} \frac{1}{2} |f_\alpha(\mathbf{x}) - y| dP(\mathbf{x}, y) \end{aligned}$$

The second integral is 0 because on the domain of integration,

$$\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) = y\}$$

the integrand, i.e., the function being integrated, is 0.

For the first integral, we note that on its domain of integration,

$$\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\},$$

its integrand is always 1. Therefore,

$$\begin{aligned} \int_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}} \frac{1}{2} |f_\alpha(\mathbf{x}) - y| dP(\mathbf{x}, y) &= \int_{\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}} 1 dP(\mathbf{x}, y) \\ &= P(\{(\mathbf{x}, y); f_\alpha(\mathbf{x}) \neq y\}) \end{aligned}$$

which, according to equation (1), is exactly $err_P(f_\alpha)$.

In the PAC approach to learning, the selection of a particular hypothesis, f_α is done such that with *high probability the risk is very small*, and this makes this learning *probably* (high probability) *almost correct* (small risk).

Furthermore, it has been customary to

1. fix the bounds on the error and probability of error and
2. to calculate the sample size, i.e. the size of the training set, **needed** to reach these bounds.

This is usually unrealistic, first, because in a learning problem we have a training set, “**AS IS**” and most often cannot get more examples, and second, because it turns out that in order to reach good results (tight bounds) the needed sample size turns out to be so large that for reasonable sample sizes (which would be available) the errors and their probabilities are too big to be useful.

More precisely, given S a set of training examples,

$$P(S; R(\alpha) = \text{err}_P(f_\alpha) \leq \epsilon) \geq 1 - \delta \quad (2)$$

which is equivalent to (use $P(\bar{A}) = 1 - P(A)$, for any event A)

$$P(S; R(\alpha) = \text{err}_P(f_\alpha) > \epsilon) < 1 - \delta \quad (3)$$

Recall that a hypothesis is said to be *consistent* if this hypothesis is correct with respect to the training set S . This means that the empirical risk, R_{emp} must be 0 and hence the bound on $R(\alpha)$ will be

$$R(a) \leq \sqrt{\frac{h(\ln \frac{2m}{h} + 1) - \ln(\frac{\epsilon}{4})}{m}} \quad (4)$$

in which, h denotes the VC¹ dimension of the hypothesis space H from which f_α comes, and m is the size of the training set S . Note then that for consistent hypotheses the bound on $R(\alpha)$ depends only on the ratio $\frac{h}{m}$.

¹Named after the Russian mathematicians Vapnik & Chervonenkis.

Let now f_α be a hypothesis consistent to a training example $(x_0, y_0) \in S \times \{-1, +1\}$ such that

$$P((x_0, y_0); R(\alpha) > \epsilon) < 1 - \epsilon \quad (5)$$

Therefore, by the iid assumption, for all $(\mathbf{x}, y) \in S \times \{-1, +1\}$

$$P((\mathbf{x}, y); R(\alpha) > \epsilon) = \Pi_{(\mathbf{x}, y)} P((\mathbf{x}, y); R(\alpha) > \epsilon) = (P((\mathbf{x}, y); R(\alpha) > \epsilon))^m < (1 - \epsilon)^m \quad (6)$$

and since $(1 - \epsilon)^m < e^{-\epsilon m}$ it follows that

$$P((\mathbf{x}, y) \in S \times \{-1, +1\}; R(\alpha) > \epsilon) < e^{-\epsilon m} \quad (7)$$

Now let us denote by $|H|$ the cardinality of the hypothesis space, that is the number of functions f_α (we assume that H is finite). Then the probability that a hypothesis f_α which satisfies (7) exists is:

$$\begin{aligned} & P\left(\bigcup_{f \in H} \{(\mathbf{x}, y) \in S; \text{ } h \text{ consistent to } S; R(\alpha) > \epsilon\}\right) \\ &= \sum_{f \in H} P(\{(\mathbf{x}, y) \in S; \text{ } h \text{ consistent to } S; R(\alpha) > \epsilon\}) \leq |H| \cdot e^{-\epsilon m} \end{aligned}$$

In order for equation (2) (or (3)) to hold we require then

$$|H|e^{-\epsilon m} = \delta \Leftrightarrow \ln(|H|e^{-\epsilon m}) = \ln \delta$$

$$\Leftrightarrow$$

$$\ln(|H|) + \ln(e^{-\epsilon m}) = \ln \delta \Leftrightarrow \ln(|H|) - \ln(\epsilon m) = \ln \delta$$

$$\Leftrightarrow$$

$$\ln(|H|) - \ln \delta = \epsilon \cdot m \Leftrightarrow \ln \frac{|H|}{\delta} = \epsilon \cdot m$$

$$\Leftrightarrow$$

$$\epsilon = \frac{1}{m} \ln \frac{|H|}{\delta} \quad (8)$$

or,

$$m = \frac{1}{\epsilon} \ln \frac{|H|}{\delta} \quad (9)$$

That is, **the sample size is proportional to the size of the hypothesis space**. Thus,

$$P \left(S; f_\alpha \text{ consistent to } S; R(\alpha) > \frac{1}{m} \ln \frac{|H|}{\delta} \right) < \delta \quad (10)$$

Let us see what this means from the point of view of PAC learning:

Assume that $\delta = 0.05$ and $\epsilon = 0.01$. Then, according to the equation (9), we need

$$m = \frac{1}{0.01} \ln \frac{|H|}{0.05} = 100 \ln(500|H|).$$

(Again, this means that the sample size is proportional to the size of the hypothesis space.)

For example,

- If $|H| = 3$ then we obtain $m \simeq 100 \times 7.31 = 731$;
- If $|H| = 10$ we obtain $m \simeq 100 \times 8.51 = 851$, etc.

Note however, that in general, $|H|$ will be much larger.

Let us now look at the following theorems:

Theorem 1 *Let H be a hypothesis space with finite VC dimension h . Then for any distribution P on $\mathbf{X} \times \{-1, +1\}$:*

$$P \left[S; \exists f \in H; R_{emp} = 0; R(\alpha) < \frac{2}{m} \left(\frac{h \ln(m) + 1}{h} + \ln \frac{2}{\epsilon} \right) \right] = 1 - \delta$$

provided that $h \leq m$, $m > \frac{2}{\epsilon}$.

Theorem 2 *Let H be fixed with finite VC dimension h . Then for every $f \in H$ there exists \mathbf{P} such that*

$$P \left[S; R_{emp} = 0; R(\alpha) \geq \max \left\{ \frac{h-1}{32m}, \frac{1}{m} \ln \left(\frac{1}{\delta} \right) \right\} \right] \geq \delta \quad (11)$$

Looking at equation 11 on the complementary we have

$$P \left[S; R_{emp} = 0; R(\alpha) < \max \left\{ \frac{h-1}{32m}, \frac{1}{m} \ln \left(\frac{1}{\delta} \right) \right\} \right] < 1 - \delta \quad (12)$$

In general, the quantity $\frac{1}{m} \ln \left(\frac{1}{\delta} \right)$ will be very big (since we want δ very small).

However, if h is big enough, more precisely, if

$\frac{h-1}{32m} > \frac{1}{m} \ln \left(\frac{1}{\delta} \right)$ (which is equivalent to $h > 1 + 32 \ln \left(\frac{1}{\delta} \right)$ which can be very large!),

then the maximum of the two quantities in the equations above is actually

$$\frac{h-1}{32m}$$

and so by setting this quantity equal to ϵ we obtain

$$m = \frac{h-1}{32\epsilon},$$

which means that for large VC dimension, to obtain small error m must be very large.

Example 1 For example, $\epsilon = 0.01$, $\delta = 0.02$ and assume $h > 1 + 32 \ln 1/\delta \simeq 126$, say $h = 127$. We need

$$m = \frac{126}{32 \cdot 0.01} = \frac{99}{0.32} = 393 \text{ (or, } m \simeq 3h \text{)}$$

Theorems 1 and 2 should be “read” together: In both the hypothesis space, that is its VC dimension h is fixed and it is finite.

Theorem 1 says that for **every \mathbf{P} (probability distribution)**, **there exists \mathbf{f} in \mathbf{H}** such that with large probability its risk is small;

Theorem 2 says that **for every \mathbf{f} in \mathbf{H} there exists a probability distribution \mathbf{P}** such that probability to have a small risk is less than 1.

So, this means that once we fix H we should look for an approach which can find a “good” distribution (this distribution would have to be in a feature space not on the initial example space where the distribution, unknown as it maybe, is fixed).

Such an approach is taken in **Support Vector Machines(SVM)**.

The idea behind the SVM approach is to fix the class of hypotheses,

$$H = \{\mathbf{w} \cdot \mathbf{x} + b; (\mathbf{w}, b)\},$$

eventually mapping the original feature space into another feature space – of much higher dimension) and find that “good” distribution (still unknown) for which 1 holds.

The idea behind learning with SVM is that the learning will find that feature space where such a distribution exists.

Let us now look at the case in which $R_{emp} \neq 0$, i.e., there are no errors on the training set.

$$R_{emp} = \frac{1}{2m} \sum_{i=1}^m |f_{\alpha}(\mathbf{x}_i) - y_i| = \frac{k}{m},$$

where k is the cardinality of the subset of the example set for which $f_{\alpha}(\mathbf{x}_i) \neq y_i$, i.e. k is the number of training errors.

Then we have

$$P(R(\alpha)) \leq \frac{k}{m} + \sqrt{\frac{h(\ln \frac{2m}{h} + 1) - \ln \frac{\epsilon}{4}}{m}} = 1 - \epsilon$$

Note that the quantity

$$\sqrt{\frac{h(\ln(\frac{2m}{h}) + 1) - \ln(\frac{\epsilon}{4})}{m}}$$

is fixed once H and ϵ are fixed, and so the only way to keep a small bound on $R(\alpha)$ is to minimize k , the number of training errors.

This is nice, isn't it, since it kind of confirms our intuition when we try to develop any kind of learning that we should minimize the training errors.

SOME EXAMPLES

Example 2 Consider the k th nearest classifier. The basic idea in this classifier is as follows:

- For each test point x we look at its k nearest neighbors which are in the

training set and assign the test point the highest frequency label which occurs in this group.

- When $k = 1$ the algorithm assigns to each point the label of its nearest training example.

Of course this requires us to exclude the case when the nearest neighbors are points of different classes.

Then **any number of points will be learned** by this algorithm, and hence the the corresponding function set has $VC - \text{dimension} = \infty$ and $R_{emp} = 0$.

This means that the bound on R_α provides no information. Yet this type of classifier can still perform well.

Example 3 The notebook classifier. This is a classifier for which the bounds should hold but which violates the bound.

To achieve these two conditions, we would want R_α to be as large as possible and the bound on it, $R_{emp} + VC - \text{confidence}$ to be as small as possible.

That is, we want a family of classifiers which gives the worst possible actual risk, and this is 0.5, the best empirical risk, $R_{emp} = 0$, up to some number m of training examples, and such that its VC dimension, $h \leq l$ (where l is the total number of training examples) is easy to compute.

The notebook classifier is as follows: assume that the notebook contains enough room to write down the classes for m training examples.

For all others examples, the classifier will assign the same class. It is assumed that there as many positive and negative examples (equal probability to draw one of them at random).

It is clear that this classifier will have :

1. $R_{emp} = 0$ for up to m examples;
 2. $R_{emp} = 0.5$ for the remaining examples;
 3. $R_\alpha = 0.5$;
- and the VC dimension,*
4. $h = m$.

From item (1) and (2) above it follows that $R_{emp} = \frac{l-m}{2 \cdot m}$

Therefore, it follows that

$$\frac{m}{4 \cdot m} \leq \ln \frac{2 \cdot l}{m} + 1 - \frac{1}{m} \ln \frac{\epsilon}{4} \quad (13)$$

Now if we take :

$$f(t) = C \exp^{\frac{t}{4}-1} \leq 1, \quad 0 \leq t \leq 1$$

since $f(t)$ is monotonic increasing and $f(1) = 0.236$ it follows that (13) holds for all ϵ .

Structural Risk Minimization (SRM)

We see that the terms that appear in the bound on the actual risk R_α depend on the chosen class of functions (the VC confidence) and on a **particular** function selected by the training step (the empirical risk).

The idea behind SRM is to find a subset of the set of functions which minimizes the empirical risk. Since the VC dimension h is an integer, the VC confidence is not a smooth function of h .

So, the hypothesis space is divided into nested subsets of functions ordered by the VC dimension: $H_1 \subset, \dots, \subset H_n$ with corresponding VC dimensions $h_1 \leq, \dots, \leq h_n$ and such that for each subset either the VC dimension or a bound on it can be calculated. For each subset, the goal of training is to minimize the empirical risk. Let $R_{emp}(H_i)$ be the minimum empirical risk for the i th subset. Then the final selection of the function f (the learning machine) is i_0 such that

$$R_{emp}(H_{i_0}) + VC(h_0) = \min_i \{R_{emp}(H_i) + VC(h_i)\}.$$

Example 4 (*Conjunction of Booleans are PAC learnable*) C : the class of target concepts – conjunction of boolean literal – (e.g. $a \wedge b \wedge \text{not } c$).

Claim 1 C is PAC-learnable

Proof:

Any consistent learner requires only a polynomial number of training examples: L consistent learner, H hypothesis space is identical to C .

Then from equation (9) we can compute m . All is needed is size of H , $|H|$ which is 3^n since each variable can (i) be included in a hypothesis, (ii) its negation can be included, or (iii) ignored.

Then from (9) with $H = 3^n$ we obtain

$$m \geq \frac{1}{\epsilon} \left(n \ln 3 + \ln \left(\frac{1}{\delta} \right) \right)$$

For example, $n = 10$, $\delta = 0.05$, $\epsilon = 0.1$), that is, with probability 0.95 the hypothesis will learn the concept with error up to 0.1, then

$$m = \frac{1}{0.1} (10 \ln 3 + \ln(1/0.05)) = 140$$

m grows linearly with n , $\frac{1}{\epsilon}$ and logarithmically in $\frac{1}{\delta}$.