

Machine Learning Theory, Self-Study

Selected Solutions for Mohri's Foundation to ML (2nd)

Others to be added on an ad hoc basis

To the curious: ML is not the same thing as deep learning. Theory here does not include neural network work and requires more human intervention by nature. The latter will be included in its own repo with PyTorch scripts. Algorithmic analysis and probabilistic considerations found in CLRS text work under dev.

Guilherme Albertini

Contents

1	PAC Learning	1
1.1	Priors	1
1.2	Exercises	1
1.3	Summary	5

Preface

TODO blah blah

Chapter 1

PAC Learning

1.1 Priors

The appendix is quite helpful for bounds used throughout the text. As the authors state, “The book of Kearns and Vazirani (1994) is an excellent reference dealing with most aspects of PAC-learning and several other foundational questions in machine learning. Our example of learning axis-aligned rectangles, also discussed in that reference, is originally due to Blumer et al. (1989).”

Mohri’s course notes will be helpful: <https://cs.nyu.edu/~mohri/ml20/>. Note that the textbook pdf is freely available on his website, too. Familiarity with data science fundamentals, multivariable calculus, linear algebra, and algorithmic analysis is assumed for this specific text. Knowledge of convex optimization (or nonlinear programming) and real analysis would be very useful, it seems.

I include a “corrected” and (imo) clearer proof from other authors that do not assume continuity of distribution under the chapters folder: http://compbio.fmph.uniba.sk/vyuka/ml/handouts/rectangles_correction.pdf

For proofs throughout: For implication $p \implies q$, an antecedent (or hypothesis) p is a sufficient condition for a consequent (a conclusion) q when the truth of p alone implies the truth of q ; however, p being false does not always imply that q is also false. A necessary condition is when the truth of q is guaranteed by the truth of p , or we can say that the truth of p is implied by the truth of q ; in other words, p is not possible without q . Several necessary conditions may induce a condition whereas a sufficient condition is alone enough to produce the said condition. The sufficient term is the part that immediately follows “if” and the necessary term is the part that immediately follows the “then”. Note that these are converses, but the converse may not always be true.

Further reading:

<https://philosophy.stackexchange.com/questions/22/what-is-the-difference-between-necessary-and-sufficient>
<https://www.kaptest.com/study/lsat/lsat-formal-logic-necessary-vs-sufficient/>
<https://pages.cs.wisc.edu/~shuchi/courses/787-F07/scribe-notes/lecture25.pdf>
https://www.cs.princeton.edu/courses/archive/spr06/cos511/scribe_notes/0214.pdf
<https://jeremykun.com/2014/01/02/probably-approximately-correct-a-formal-theory-of-learning/>
<https://www.cs.cornell.edu/courses/cs6781/2020sp/lectures/03-pac1.pdf>

1.2 Exercises

Exercise 2.1

Necessary and sufficient: a concept class \mathcal{C} is efficiently PAC-learnable using hypothesis space \mathcal{H} in the standard PAC model if and only if it is efficiently PAC-learnable using the hypothesis space $\mathcal{H} \cup \{h_0, h_1\}$ in the two-oracle PAC model.

Sufficiency: Show that if \mathcal{H} is PAC-learnable in the standard, one-oracle model then so too is it in the variant.

Since \mathcal{C} is efficiently PAC-learnable using \mathcal{H} , there exists an algorithm \mathcal{A} and a polynomial p such that for any distribution \mathcal{D} and any target concept $c \in \mathcal{C}$, if \mathcal{A} is given a sample of size $m \geq p(1/\epsilon, 1/\delta, n, \text{size}(c))$ drawn from \mathcal{D} , it outputs a hypothesis h from \mathcal{H} such that with probability at least $1 - \delta$, the error of h on \mathcal{D} is at most ϵ . Assume a distribution \mathcal{D} on $\mathcal{X} \times \{-1, +1\}$. The learner's goal is to output a hypothesis with such probability over the choice of two training sets (in the stochastic scenario where the output label is a probabilistic function of the input, thus does not guarantee unique labels) and requires both $\mathbb{P}[R(h)_{x \sim \mathcal{D}^+} \leq \epsilon] \geq 1 - \delta$ and $\mathbb{P}[R(h)_{x \sim \mathcal{D}^-} \leq \epsilon] \geq 1 - \delta$. By the law of total probability for the error rate (or risk),

$$\begin{aligned} R(h)_{\mathcal{D}^m} &= \mathbb{E}_{x \sim \mathcal{D}^m} [\mathbb{1}_{h(x) \neq c(x)}] = \sum_{x \in \mathcal{D}^m} x \mathbb{P}[\mathbb{1}_{h(x) \neq c(x)}] = \mathbb{P}_{x \sim \mathcal{D}^m} [h(x) \neq c(x)] \\ &= \mathcal{D}^m(\mathcal{X}^+)(R(h)_{\mathcal{D}^+}) + \mathcal{D}^m(\mathcal{X}^-)(R(h)_{\mathcal{D}^-}) \end{aligned}$$

Now for a weighted sampling method from the negative and positive instance distributions we can say $\epsilon_{\mathcal{D}} = \mathbb{P}[\mathcal{D}^+](\alpha\epsilon)_{\mathcal{D}^+} + \mathbb{P}[\mathcal{D}^-](\beta\epsilon)_{\mathcal{D}^-} = \epsilon \underbrace{(\mathbb{P}[\mathcal{D}^+](\alpha - \beta) + \beta)}_{\leq 1}$ for some constants $0 \leq \alpha, \beta \leq 1$. Note that

$\alpha = \beta = 0 \implies \epsilon = 0$, which means you're demanding that the learned hypothesis perfectly match the true target function, which can lead to overfitting and complex hypotheses that are computationally expensive to work with (recall $\epsilon > 0$ by definition). Conversely, by setting $\epsilon = 1$, you're ok accepting any hypothesis without regard to its performance. Thus, select an appropriate δ and, for simplicity, set $\mathbb{P}[\mathcal{D}^+] = \frac{1}{2}$, to get

$$\begin{aligned} \mathbb{P}[R(h)_{\mathcal{D}} \leq \frac{\alpha\epsilon}{2}] &\geq 1 - \delta \\ \frac{1}{2}(\mathbb{P}[R(h)_{\mathcal{D}^+} \leq \alpha\epsilon] + \mathbb{P}[R(h)_{\mathcal{D}^-} \leq \alpha\epsilon]) &\geq 1 - \delta \end{aligned}$$

If we set $\alpha = 1$ for this case we then see that by selecting $\mathbb{P}[R(h)_{\mathcal{D}} \leq \frac{\epsilon}{2}]$ with an appropriate confidence interval, we must have that both $\mathbb{P}[R(h)_{\mathcal{D}^-} \leq \epsilon], \mathbb{P}[R(h)_{\mathcal{D}^+} \leq \epsilon] \geq 1 - \delta$. We immediately notice that a biased dataset would then require us to make considerable adjustments to the overall error rate, as \mathcal{D}^+ will shift the weight of distribution of samples. Another thing to note is that by setting α or β to 0, we're essentially requiring that the hypothesis have zero error on positive (or negative) instances (perhaps requiring an absurdly complex hypothesis to do so). This transforms the problem into a rather stringent one-oracle PAC model focused solely on the positive (or negative) class which is not sensitive to noise. The PAC framework is designed to find a balance between minimizing errors on both classes while accounting for uncertainties in real-world data; this decoupled mode tells the learner to query the oracle for instances of one class while imposing stringent error requirements on the other class.

Necessary: Show that if \mathcal{H} is PAC-learnable in the two-oracle variant, it is also PAC-learnable in the standard model.

We now can assume that \mathcal{C} is efficiently PAC-learnable in the two-oracle PAC model so there exists an algorithm \mathcal{A} such that for $c \in \mathcal{C}$, $\epsilon, \delta > 0$, there are m^+ and m^- in $p(1/\epsilon, 1/\delta, \text{size}(c))$, such that if we draw at least this number of negative and positive instances with confidence of at least $1 - \delta$, the hypothesis h output by the learner satisfies:

$$\begin{aligned} R(h)_{\mathcal{D}^{+,-}} &\leq \epsilon \\ \mathbb{P}[R(h)_{\mathcal{D}^{+,-}} \leq \epsilon] &\leq \mathbb{P}[\epsilon] = \epsilon \end{aligned}$$

So, given sufficient numbers of negative and positive examples, we can generate a hypothesis h such that it has low errors on both negative and positive instances. If we draw too few examples, the conclusions about

the hypothesis's performance might not hold true for the entire distribution (generalize); to bridge the gap between the variant and standard model it is then best to take $m \geq \max\{m^+, m^-\}$ and, drawing such with polynomial conditions above and using the union bound and total probability,

$$\begin{aligned}\mathbb{P}[R(h)_{\mathcal{D}}] &\leq \mathbb{P}[\mathcal{D}^+] \mathbb{P}(R(h)_{\mathcal{D}^+}) + \mathbb{P}[\mathcal{D}^-] \mathbb{P}(R(h)_{\mathcal{D}^-}) = \mathbb{P}[R(h)_{\mathcal{D}} | c(x) \neq -1] \mathbb{P}[c(x) \neq -1] + \mathbb{P}[R(h)_{\mathcal{D}} | c(x) \neq 1] \mathbb{P}[c(x) \neq 1] \\ &\leq \epsilon (\mathbb{P}[c(x) \neq -1] + \mathbb{P}[c(x) \neq 1]) \leq \epsilon\end{aligned}$$

Let X be the total number of positive examples obtained from drawing m examples, with probability of a positive example of ϵ as shown above. Note $\mathbb{E}[X] = m\epsilon$ and if $\epsilon < 1$, then $\mathbb{E}[X] < m$, indicating that you expect to obtain fewer positive examples on average than the total number of examples drawn.

$$\mathbb{P}\left[\frac{X}{m} \leq (1 - \gamma)\epsilon\right] \leq e^{-\frac{m\epsilon\gamma^2}{2}} \implies \mathbb{P}\left[X > (1 - \gamma)m\epsilon\right] \leq e^{-\frac{m\epsilon\gamma^2}{2}}$$

Setting (a rather lax) $\gamma = \frac{1}{2}$, required that $\gamma \in [0, 1/\epsilon - 1]$, to match conditions above and needing $m^+ = m(1 - \gamma)\epsilon = \frac{m\epsilon}{2}$,

$$\mathbb{P}\left[X > m^+\right] \leq e^{-\frac{m^+}{4}} \leq \frac{\delta}{2}$$

Since we want to include the minimum number of positive instances, we set the latter expression to an appropriate bound $\delta/2$ (from above, we saw that using an even split weight between negative and positive examples with overall $\epsilon/2$ as we set here, so can use $2(1 - (\delta/2)) > 1 - \delta$) and substitute back to get an expression of $m \geq \min\{\frac{2m^+}{\epsilon}, \frac{8}{\epsilon} \log(2/(\delta))\}$. A similar procedure is done with the negative case to arrive at $m \geq \min\{\frac{2m^+}{\epsilon}, \frac{2m^-}{\epsilon}, \frac{8}{\epsilon} \log(2/(\delta))\}$ for a balanced dataset.

Exercise 2.2

$$\begin{aligned}\mathbb{P}_{S \sim \mathcal{D}^m}[R(h_S) > \epsilon] &\leq \mathbb{P}_{S \sim \mathcal{D}^m}[\cup_{i=1}^{2n} \{h_S \wedge r_i = \emptyset\}] \\ &\leq \sum_i \mathbb{P}_{S \sim \mathcal{D}^m}[\{h_S \wedge r_i = \emptyset\}] && \text{(union bound)} \\ &\leq 2n(1 - \frac{\epsilon}{2n})^m && \mathbb{P}[r_i] \geq \frac{\epsilon}{2n} \\ &\leq (2n)e^{-\frac{m\epsilon}{2n}} \leq \delta \iff m \geq \frac{2n}{\epsilon} \log(2n/\delta)\end{aligned}$$

Exercise 2.3

Define annulus $A := \{a \leq \|x - x_0\| \leq r\}$ with $a := \sup\{r' \mid \mathbb{P}[r' \leq \|x - x_0\| \leq r] > \epsilon\}$. With this construction $\mathbb{P}[A] \geq \epsilon$. If the inner circle \mathcal{C}' intersects the whole annular region for a sample, then the error (disagreement between target and hypothesis) is at most ϵ . Using the contrapositive,

$$\begin{aligned}\mathbb{P}_{S \sim \mathcal{D}^m}[\text{error}_D(\mathcal{C}') > \epsilon] &\leq \mathbb{P}_{S \sim \mathcal{D}^m}[\{\mathcal{C}' \wedge A = \emptyset\} \text{ for some } i] \\ &\leq \sum_i \mathbb{P}_{S \sim \mathcal{D}^m}[\{\mathcal{C}' \wedge A = \emptyset\}] \\ &= \sum_i \mathbb{P}_{S \sim \mathcal{D}^m}[\{S \wedge A = \emptyset\}] \\ &= \sum_i (1 - \mathbb{P}[A])^m \\ &\leq (1 - \epsilon)^m \leq \delta\end{aligned}$$

and the result follows.

Exercise 2.4

The thinking applied previously suggested that the true (generalization) error being at most a small ϵ was implied by the hypothesis (inner circle) intersecting the error region. The contrapositive then implies that if the error is greater than ϵ , then the hypothesis misses part of the error region. If we just look at the three points provided as the sample, we see that the error is most likely greater than ϵ even when the hypothesis intersects the error region for all instances. This invalidates the first part of the logic so the contrapositive cannot hold to generate the PAC-learning proof.

Exercise 2.5

Assume $\mathbb{P}[ABC] > \epsilon$. We can show that the triangle formed by $A'B'C'$ contains triangle $A''B''C''$. Note that the error is at most ϵ when the area intersects all three regions defined by $A''B''C''$'s segments. As $A'B'C'$ does just this as we see with previous considerations, then $R(h_S = A'B'C')_{\mathcal{D}} \leq \epsilon$. Using the contrapositive, there must be an error greater than ϵ should we miss at least one of the three regions for which $\mathbb{P}[r_i] \geq \frac{\epsilon}{3}$.

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m}[\text{error}_D(h_S) > \epsilon] &\leq \mathbb{P}_{S \sim \mathcal{D}^m}[\{h_S, r_i\} = \emptyset] \\ &\leq \sum_{i=1}^3 \mathbb{P}_{S \sim \mathcal{D}^m}[\{h_S, r_i\} = \emptyset] \leq 3(1 - \frac{\epsilon}{3})^m \leq 3e^{-\frac{m\epsilon}{3}} \leq \delta \end{aligned}$$

And the result follows.

Exercise 2.6

The probability that each x_i in the sample will miss region r_j implies that the sample already lied outside the region and so was not flipped to negative or that the sample instance initially laid inside the region but was later flipped to negative.

$$\begin{aligned} \mathbb{P}[\{x_i \mid x_i \notin r_j \cup (x_i \in r_j, c(x_i) * 0, c(x_i) = 1)\}] &= \mathbb{P}[\{x_i \mid x_i \notin r_j \cup (x_i \in r_j, 1 * 0)\}] \\ &= \mathbb{P}[x_i \notin r_j] + \eta \mathbb{P}[x_i \in r_j] \\ &= \mathbb{P}[x_i \notin r_j] + \eta(1 - \mathbb{P}[x_i \notin r_j]) \\ &= \mathbb{P}[x_i \notin r_j](1 - \eta) + \eta \\ &\leq (1 - \frac{\epsilon}{4})(1 - \eta) + \eta \\ &\leq 1 - \frac{(1 - \eta)\epsilon}{4} \\ &\leq 1 - \frac{(1 - \eta')\epsilon}{4} \end{aligned}$$

Now use the usual PAC learning logic,

$$\mathbb{P}[R(R') > \epsilon] \leq \sum_i 1 - \frac{(1 - \eta')\epsilon}{4} \leq 4(1 - \frac{(1 - \eta')\epsilon}{4}) \leq 4e^{-\frac{m\epsilon(1 - \eta')}{4}} \leq \delta$$

So for with probability at least $1 - \delta$, we'd arrive at an approximately correct hypothesis given $m \geq \frac{4}{(1 - \eta')\epsilon} \log \frac{4}{\delta}$. This decays to the original if we remove noise.

Exercise 2.7

- For any hypothesis, let $d(h)$ denote the probability that the label of a training point received by the learner disagrees with the one given by h . The probability that the label of a point be incorrect is η by definition.

- Disagreement between the target function h^* and the hypothesis occurs when the algo fails when the label is correct and when it correctly predicts an incorrect label.

$$\begin{aligned}\mathbb{P}[\text{incorrect}] &= \mathbb{P}[c(x) \text{ correct}, c(x) \neq h(x)] + \mathbb{P}[c(x) \text{ incorrect}, c(x) = h(x)] \\ &= R(h_S)(1 - \eta) + \eta(1 - R(h_S)) = \eta + (1 - 2\eta)R(h_S)\end{aligned}$$

- Assuming that $R(h_S) > \epsilon$,

$$d(h) - d(h^*) = (1 - 2\eta)R(h_S) \geq (1 - 2\eta)\epsilon \geq (1 - 2\eta')\epsilon$$

- We use Hoeffding's inequality. Recall that $d(h^*) = \eta$ is a random variable and so $\mathbb{E}[d(h^*)] = \sum_i x_i \eta_i = m\eta = \hat{d}(h^*)$, as defined in the problem.

$$\begin{aligned}\mathbb{P}[\hat{d}(h^*) - d(h^*) \geq \epsilon'/2] &\leq e^{-\frac{2(\frac{\epsilon'}{2})^2}{\sum_{i=1}^m 1^2}} = e^{-\frac{(\epsilon')^2}{2m}} \leq \frac{\delta}{2} \\ \frac{1}{m} &\geq \frac{2}{(\epsilon')^2} \log \frac{2}{\delta} \implies m \geq \frac{2}{(\epsilon')^2} \log \frac{2}{\delta} \quad (m \text{ is positive integer})\end{aligned}$$

- Similar to before, but using union bounding.

$$\begin{aligned}\mathbb{P}[h \in \mathcal{H} \mid d(h) - \hat{d}(h) \geq \epsilon'/2] &\leq |\mathcal{H}| e^{-\frac{2(\frac{\epsilon'}{2})^2}{\sum_{i=1}^m 1^2}} = |\mathcal{H}| e^{-\frac{(\epsilon')^2}{2m}} \leq \frac{\delta}{2} \\ \frac{1}{m} &\geq \frac{2}{(\epsilon')^2} (\log \frac{2}{\delta} + \log |\mathcal{H}|) \implies m \geq \frac{2}{(\epsilon')^2} (\log \frac{2}{\delta} + \log |\mathcal{H}|) \quad (m \text{ is positive integer})\end{aligned}$$

- Using the bounding hint provided,

$$\begin{aligned}\hat{d}(h) - \hat{d}(h^*) &= \underbrace{[\hat{d}(h) - d(h)]}_{\geq -\frac{\epsilon'}{2}} + \underbrace{[d(h) - d(h^*)]}_{\geq \epsilon'} + \underbrace{[d(h^*) - \hat{d}(h^*)]}_{\geq -\frac{\epsilon'}{2}} \\ &\geq -\frac{\epsilon'}{2} + \epsilon' - \frac{\epsilon'}{2} = 0\end{aligned}$$

The fraction of the points in S whose labels disagree with those given by any $h \in \mathcal{H}$ is at least (likely greater) than that for the target function. Thus, the learner will not select this set of poor hypotheses as it cannot minimize this difference bound. We note this learner can be used for PAC-learning despite noise in the consistent (“realizable”) case where we make the strong assumption the target $c \in \mathcal{H} \subset \mathcal{C}$.

1.3 Summary

PAC learning makes a rather strong assumption that there exists a function (concept), which is a subset of the domain, $c : \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathcal{Y} = \{0, 1\}$. We assume that examples are independently and identically distributed (i.i.d.) according to some fixed but unknown data distribution and the learner considers a set of concepts (hypotheses, set \mathcal{H}) which may or may not coincide with the “true” set of concepts (\mathcal{C}). It receives a sample S drawn i.i.d. according to \mathcal{D} as well as the labels $(c(x_1), \dots, c(x_m))$, which are based on a specific target concept $c \in \mathcal{C}$ to learn. The task is then to use the labeled sample S to select a hypothesis $h \in \mathcal{H}$ that has a small generalization error with respect to the concept c .

While it would be great to always select a hypothesis for which we get 0 error, since we draw from a training set which have a non-zero probability of have misleading examples (thus bound such cases with δ), there may be several hypothesis consistent with such a sampling; the only way for the learner to pick the optimal solution for the target is to train on all the samples in the domain, or universe (unrealistic). It may also just overfit to memorize instances in a sampling to become a baseline classifier. We relax this assumption with the small ϵ . Note that $\epsilon = 0$ implies that the the learned hypothesis needs to perfectly match the true target function, which can lead to this overfitting and possibly require complex hypotheses.

The converse $\epsilon = 1$ would allow for any hypothesis without regard to its performance, hence no learning. To further relax the assumption that (near) zero error be made by all predictors in \mathcal{H} (without knowing if the target concept c lies in \mathcal{H}), we insist on a finite set of hypotheses, which requires some inductive bias; that is, we make an assumption that there exists a “good” set of candidate predictors which contain the target concept, $c \in \mathcal{H}$. As before, PAC learning then guarantees an approximately correct (small true loss) hypothesis (from a possible set of minimizers) given a sufficiently sized sample. In the most (realistic) general case, there may be no hypothesis in \mathcal{H} consistent with the labeled training sample. This gives us more latitude for a fixed $|\mathcal{H}|$, but to attain the same guarantee as in the consistent case, a quadratically larger labeled sample is needed in this tradeoff between error and sample size (see Mohri).

$$\begin{aligned}
& \mathbb{P}_{S \sim \mathcal{D}^m} [h > \epsilon \text{ inaccurate on average}] \stackrel{\Longleftrightarrow}{\substack{\text{(in-sample error 0, but not for universe)}}} \mathbb{P}_{S \sim \mathcal{D}^m} [\text{algo fooled}] \\
& \mathcal{H}_{\text{bad}} = \{\exists h \in \mathcal{C} \mid h(x) = c(x) \ \forall x \in S, \text{ but } \mathbb{P}_{S \sim \mathcal{D}^m} [h(x) \neq c(x)] \geq \epsilon\} \\
& \mathcal{H}_{\text{good}} = \{\exists h \in \mathcal{C} \mid h(x) = c(x) \ \forall x \in S, \text{ and } \mathbb{P}_{S \sim \mathcal{D}^m} [h(x) \neq c(x)] \leq 1 - \epsilon\} \\
& \mathbb{P}_{S \sim \mathcal{D}^m} [\mathcal{H}_{\text{bad}} \text{ non-empty, so algo may output one of its members}] \leq \delta
\end{aligned}$$

For given S , $\mathcal{H}_{\text{good}}$ is never empty as it must at least contain the target concept (later relaxed, only for “realizable” version)

Note that the proof of consistent, finite set case relied on the underlying technique of sampling from the data distribution and then determining a good generalizable hypothesis after assessing all in-sample errors. That is, there may be a set of (poor) samples (x_i ’s) for which the true risk under a chosen hypothesis is larger than ϵ yet have 0 empirical risk for a sample. We would want to avoid the set of all bad hypotheses for which this “bad” condition holds. The set of all samples that will lead the chosen (bad) hypothesis to produce error greater than ϵ is a subset of the set of all samples for which there exists some (bad) hypothesis that produces zero in-sample error (which we may or may not have chosen). Thus, by bounding the superset (and thus further constricting the subset of interest), we restrict such hypotheses becoming godly predictors for all instances in a sample, thus bounded by a small δ .

Wrapping up: good learners will learn with high probability and close approximation to the target concept. With high probability, we find a hypothesis that will have low error (approximately correct), requiring parameters $\epsilon, \delta > 0$. Thus, with probability of at least $1 - \delta$, the algorithm learns the concept with error at most ϵ . The ϵ is the upper bound on the accuracy ($1 - \epsilon$). We use δ to bound the probability of failure in achieving the set accuracy (a bad event) and thus want a hypothesis to be approximately correct with confidence $1 - \delta$. PAC is powerful in that it requires no assumption about the underlying distribution but has the rather unrealistic requirement that the labeling rule (target concept) is in the hypothesis space. Agnostic (i.e., not realizable) PAC learning relaxes this assumption; in scenarios for which there may exist no such concept, and instead that all distributions produce labels stochastically (and thus may not be unique), we are left with a minimal non-zero error for any hypothesis.

