

CS 2280: Computational Learning Theory
Homework 1
Due: Feb. 19, 11:59pm
Nico Fidalgo

Note For some problems you may need the following fact from probability theory called the Chernoff bound:

Theorem 1 (Chernoff bound). *Let X_i , $1 \leq i \leq m$, be independent Bernoulli random variables each with probability of success μ . Then for any $\lambda \in [0, 1]$,*

$$\Pr \left[\left| \frac{1}{m} \sum_{i=1}^m X_i - \mu \right| \geq \lambda \mu \right] \leq 2e^{-\lambda^2 \mu m / 2}$$

You can find more details and generalizations in the textbook or online.

Policy reminders You are strongly encouraged to type your solutions using L^AT_EX. You may discuss problems with your classmates, but not merely copy each others solutions. You must write all solutions by yourself, list your collaborators on your problem sets and also appropriately cite any resources outside of the class materials that you have used. You are not allowed to look up solutions to the problems. Please do not use LLMs or LLM-assisted tools for finding solutions to the problems.

Problem 1. (10pt) **Learning union of intervals.** Give a PAC-learning algorithm for the concept class \mathcal{C}_2 formed by unions of two closed intervals, that is $[a, b] \cup [c, d]$, with $a, b, c, d \in \mathbb{R}$. Extend your result to derive a PAC-learning algorithm for the concept class \mathcal{C}_n formed by unions of $n \geq 1$ closed intervals, thus $[a_1, b_1] \cup \dots \cup [a_n, b_n]$, with $a_i, b_i \in \mathbb{R}$ for $k \in [n]$. What are the time and sample complexities of your algorithm as a function of n ?

To give a PAC-learning algorithm for the union of n intervals, the goal is to learn a function that correctly classifies points as either inside or outside the target intervals while ensuring an error of at most ϵ with confidence at least $1 - \delta$.

Given a set of positive samples $S = \{x_1, x_2, \dots, x_m\}$, we construct the hypothesis by identifying pairs of adjacent positive samples. We assume there are negative samples in between the intervals to separate them, $\{c_0, c_1, \dots, c_n\}$. The hypothesis $h(x)$ labels x as 1 if it lies in an interval $[x_i, x_j]$ and 0 otherwise.

To control the error, we must ensure that the positive samples sufficiently cover the intervals, particularly near their boundaries. If an interval lacks positive samples near its edges (in a fraction of length $\frac{\epsilon}{2}$), it may lead to classification errors. The probability of missing samples in any of these $2n$ regions is bounded by

$$P(\text{missing any of } 2n \text{ tails}) \leq 2nP(\text{missing one tail})$$

Which, using probability bounds, reduces to

$$P(\text{miss one tail}) \leq e^{-\frac{m\epsilon}{2}}$$

At most, this probability is δ , therefore the sampling complexity is

$$\frac{2}{\epsilon} \ln \frac{2n}{\delta} \leq m$$

This means that the number of positive samples needed grows logarithmically with the number of intervals n , i.e., $O(\log n)$.

The number of negative samples needed is at most $O(n)$ since we only need to place one between each pair of adjacent intervals. The time complexity of processing all samples is therefore $O(m + n)$ which simplifies to $O(n)$ as m grows as $O(\log n)$.

Problem 2. (10pt) **PAC-learning parity functions.** *Parity functions* is analogous to conjunctions or disjunctions, but with the XOR operation \oplus rather than \wedge or \vee . To be precise, given a subset $S \subseteq \{1, 2, \dots, n\}$, a parity function $f_S : \{0, 1\}^n \rightarrow \{0, 1\}$ outputs 1 iff the number of 1s in the subset of the input indexed by S is odd. For example, if $n = 4$ and $S = \{1, 3, 4\}$, the parity function $f_S = x_1 \oplus x_3 \oplus x_4$ will return $f(0010) = 1, f(1010) = 0$. The class of parity functions \mathcal{P} includes all functions that can be described in this way:

$$\mathcal{P} = \{f_S : S \in 2^{\{1, \dots, n\}}\} \quad f_S(x) = \bigoplus_{i \in S} x_i$$

Show that \mathcal{P} is PAC learnable.

Hint. Note that $x \oplus y$ is equal to $x + y \pmod{2}$. Additionally, searching for an Occam algorithm might be easier than working with ϵ and δ .

Each parity function is defined by a subset $S \subseteq \{1, \dots, n\}$. Learning the function means discovering the unknown set S using examples $(x, f_S(x))$ where x is an n -bit binary vector and $f_S(x)$ is the parity (XOR sum) of the bits in S , i.e., $f_S(x) = \bigoplus_{i \in S} x_i$. The goal is to figure out S from training data.

The XOR operation is equivalent to addition modulo 2, i.e.,

$$x \oplus y = (x + y) \pmod{2}$$

If we collect enough labeled examples $(x, f_S(x))$ then we can solve a system of linear equations to determine S . To learn S , we observe that each training example $(x, f_S(x))$ gives us one linear equation with n unknowns, i.e.,

$$x_{i_1} \oplus x_{i_2} \oplus \dots \oplus x_{i_k} = f_S(x)$$

The algorithm is to collect n independent samples $(x, f_S(x))$ and use them to construct a system of n linear equations in n variables and solve for S using Gaussian elimination.

Solving for S requires at most n independent equations, thus we only need $O(n)$ samples. Gaussian elimination runs in $O(n^3)$ time, so the overall time complexity is $O(n^3)$. As Gaussian elimination is deterministic, solving for S is guaranteed to be correct, therefore satisfying the PAC-learning requirement.

Problem 3. (10pt) **“Relaxed” δ -dependence in PAC definition.** Prove that if a concept class \mathcal{C} is PAC-learnable by class \mathcal{H} by the usual definition (in time polynomial in $1/\delta, 1/\epsilon, n, \text{size}(c)$), then it is in fact learnable by class \mathcal{H} in time polynomial in $\log(1/\delta), 1/\epsilon, n, \text{size}(c)$. (Notice the more stringent requirement on the running time’s dependence on δ .)

Hint. Can we boost our confidence by running the original PAC-learning algorithm multiple times?

We start with a PAC-learning algorithm that outputs a hypothesis function h satisfying the PAC-learning requirements. Instead of relying on a single hypothesis, we run the algorithm k times and gather multiple candidate hypotheses h_1, h_2, \dots, h_k . Each hypothesis is independently trained on different data samples.

To create a more robust classifier, we can use majority voting. Given an input x , we compute $h_i(x)$ for each of the k hypotheses and the final output is determined by whichever label appears most frequently. Since each hypothesis is independent, the probability of many of them making an error at the same time follows a Chernoff bound.

Using Chernoff bounds, we ensure that the probability of the majority hypothesis failing decreases exponentially with k , leading to a failure probability at most δ when

$$k = O\left(\frac{1}{\epsilon} \log(1/\delta)\right)$$

By running the PAC-learning algorithm k times, we reduce the confidence dependence from $1/\delta$ to $\log(1/\delta)$, ensuring that the overall runtime remains polynomial in $1/\delta, 1/\epsilon, n, \text{size}(c)$.

Problem 4. (10pt) **Approximate Occam algorithm.** For constants $\alpha \geq 1$, $\beta < 1$, and $\gamma < 1/2$, let a γ -approximate (α, β) -Occam algorithm be the same as an (α, β) -Occam algorithm, but instead of consistency with all the examples only consistency with at least $1 - \gamma$ fraction of examples is guaranteed. Similarly let a γ -PAC learning algorithm for a concept class \mathcal{C} be the same as a PAC learning algorithm except that it produces a $(\gamma + \epsilon)$ -approximate (instead of an ϵ -approximate) hypothesis.

Show that if there is a γ -approximate (α, β) Occam algorithm for \mathcal{C} then there is a γ -PAC learning algorithm for \mathcal{C} .

We are given a γ -approximate Occam algorithm, meaning it finds a hypothesis h that is not perfectly consistent but misclassifies at most a γ fraction of the training examples and the hypothesis space H is restricted in size. The goal is to show that with high probability the Occam algorithm does not select a bad hypothesis when trained on a finite sample of size m , still allowing for PAC-learning, meaning we can find a hypothesis with error at most $\gamma + \epsilon$ with high probability.

A hypothesis is bad if the true error $P(h(x) \neq f(x))$ exceeds $\gamma + \epsilon$. If we sample enough data points, we should be able to identify and eliminate bad hypotheses with high probability. The probability of a bad hypothesis looking good on a small sample is exponentially small via Chernoff bounds.

The Occam algorithm restricts the number of hypotheses $|H|$. We apply the union bound to make sure no bad hypothesis is chosen, ensuring that

$$P(\text{choosing a bad } h) \leq |H|e^{-2m(\epsilon-\gamma)^2}$$

which is at most δ . Solving for m , we get

$$O\left(\frac{1}{(\epsilon - \gamma)^2} \log \frac{|H|}{\delta}\right) = m$$

By substituting the Occam complexity bound for $|H|$, we ensure the sample complexity remains polynomial in $1/\epsilon$ and $\log(1/\delta)$, satisfying the PAC-learning requirement.

Problem 5. (15pt) **Learning unions of halfspaces is NP-hard.** A *linear threshold function* or a *halfspace* is a function $f(x) : \{0, 1\}^n \rightarrow \{0, 1\}$ such that there exists some n -dimensional real vector w and real number θ such that $f(x) = 1$ iff $w \cdot x \geq \theta$. A union of halfspaces is simply a

disjunction of several halfspaces. Suppose that for some $k \geq 3$, unions of k halfspaces are efficiently PAC learnable by an algorithm that outputs unions of k halfspaces. Show that then there is a randomized algorithm R taking as input a graph G that runs in time polynomial in the size of G with the following properties:

- If G is k -colorable then with probability at least $2/3$, R outputs a k -coloring of G .
- If G is not k -colorable, then R 's output may be arbitrary.

Hint. Consider adding the example $\langle 1^n, - \rangle$ to the reduction you saw in class.

The goal is build a randomized polynomial-time algorithm R which, on an input graph G , outputs a k -coloring with high probability if G is k -colorable and something potentially arbitrary if G is not k -colorable. To do this, we assume there is a polynomial-time PAC learner L that learns unions of k halfspaces with high accuracy.

Let $G = (V, E)$ have n vertices. We represent a coloring of G by a binary vector $x \in \{0, 1\}^m$. Positive examples correspond to valid k -colorings (no adjacent vertices share the same color), and negative examples correspond to invalid colorings (at least one edge (u, v) has both endpoints colored the same). If G is k -colorable, one can build a disjunction of k linear threshold functions that accept the valid colorings and reject the invalid ones. The claim is that there exists a union of k halfspaces that separates valid from invalid colorings iff G is k -colorable.

Algorithm R picks random assignments x of colors to vertices from some easy-to-sample distribution. We label each x as $+$ if it is a valid coloring and $-$ otherwise. To avoid trivial classifiers, we add a special negative example $\langle 1^m, - \rangle$ so that any hypothesis that always outputs $+$ will get at least one error. We draw polynomially many labeled examples and feed them to the PAC learner L .

If G is k -colorable, then there is a perfect union of k halfspaces concept that classifies all training examples correctly. By PAC guarantees, with probability at least $2/3$ the learner L will find a hypothesis h close to that perfect concept. Because h classifies valid colorings positively and invalid ones negatively, we can query h on particular color assignments to recover a valid k -coloring of G . If G is not k -colorable, then there is no union of k halfspaces that is consistent with all of the examples, so L might output an arbitrary hypothesis. Thus, R builds the labeled dataset from random color assignments of G , runs L on it to obtain h , and queries h to extract a valid coloring if one exists. All steps of this algorithm can be done in polynomial time in $|G|$. Since k -coloring is NP-hard, the only way R could run in polynomial time is if the assumption that unions of k halfspaces is PAC-learnable is true, but this would imply $P = NP$, thus, learning unions of k halfspaces must be NP-hard.

Problem 6. (15pt) **Conjunction of PAC learnable classes is PAC learnable.** Show that if both class C_1 and class C_2 are PAC learnable from positive examples only (each separately), then class $C = \{c_1 \wedge c_2 \mid c_1 \in C_1 \text{ and } c_2 \in C_2\}$ is also PAC learnable from positive examples only.

Hint. Let $h = h_1 \wedge h_2$, where h_1, h_2 are obtained from feeding the positive examples of $c_1 \wedge c_2$. What do we know about the relationship between c_1 and h_1 ?

We have two separate learners L_1 and L_2 , each can learn concepts from C_1 or C_2 respectively using only positive examples from the target concept. The goal is to learn the conjunction $c_1 \wedge c_2$, i.e., a point is labeled positive iff it is positive for both c_1 and c_2 .

The first step is draw independently and identically distributed positive samples from the target concept $c = c_1 \wedge c_2$ which are guaranteed to be positive for both c_1 and c_2 . We then feed these positive samples to L_1 and L_2 which treat them as samples coming from c_1 and c_2 respectively as they each label these samples positively. Let h_1 and h_2 be the hypotheses output by L_1 and L_2 respectively. Define $h = h_1 \wedge h_2$ be the combined hypothesis that captures the intersection of the

sets recognized by h_1 and h_2 .

If h_1 is $\epsilon/2$ -accurate for c_1 on positive examples and h_2 is $\epsilon/2$ -accurate for c_2 on positive examples, then $h_1 \wedge h_2$ will be at most ϵ -inaccurate for $c_1 \wedge c_2$. The union bound handles confidence, i.e., if each learner individually succeeds with high probability, then the combined event also has high probability. The sample size and time complexity remain polynomial because you draw only one batch of positive examples for c . Both learners run in time polynomial in $1/\epsilon, 1/\delta, n, \text{size}(c)$, so the combined learner also runs in time polynomial in $1/\epsilon, 1/\delta, n, \text{size}(c)$, satisfying the PAC-learning requirement and hence, $C_1 \wedge C_2$ is PAC-learnable from positive examples only.

Problem 7. (15pt) Probabilistic PAC is equivalent to deterministic PAC. Let \mathcal{C} be any concept class and \mathcal{H} be any hypothesis class. Let h_0 and h_1 be representations of the identically 0 and identically 1 functions respectively. Show that if there is a randomized algorithm for efficiently PAC learning \mathcal{C} using \mathcal{H} , then there is a deterministic algorithm for efficiently PAC learning \mathcal{C} using $\mathcal{H} \cup \{h_0, h_1\}$.

Note: There are two sources of randomness we are talking about.

1. The inherent randomness that occurs because the example oracle EX is random.
2. The randomization that the algorithm uses. You are asked to get rid of the second randomization. More precisely if you run the learning algorithm several times using the same sequence of examples as output by the example oracle, the final hypothesis your output should be the same (and almost accurate with high probability).

This problem is complex, so we will help you organize the solution as below. Observe that a randomized algorithm can be viewed as a deterministic algorithm which takes as additional input a random string, which we may regard as a sequence of fair coin toss. The algorithm is then required to work correctly with high probability over this input string. Let the randomized algorithm be L and consider the following deterministic algorithm:

Using an oracle $EX(c, \mathcal{D})$ get M data points. Set $i = 1$.
 If less than $\epsilon M/2$ examples have label 0 return h_1 .
 Else if less than $\epsilon M/2$ examples have label 1 return h_0 .
 Else

- Whenever L requests an example return (x_i, y_i) and increment i .
- Whenever L requests a random bit as follows:
 - If $y_i \neq y_{i+1}$, then return y_i and increment i twice.
 - Otherwise, increment i twice, effectively discarding these two points.

If L terminates, return the hypothesis output by L . Else whenever $i > M$, i.e. we run out of data points, output h_0 .

There are now three sources of error from this deterministic algorithm:

- Incorrectly choosing h_0 or h_1 as the hypothesis at the start;
- L terminates but returns a hypothesis with a large error;
- L fails to terminate before running out of M examples.

For simplicity, you may assume that L requires m examples and t random bits to output with probability at least $1 - \delta/3$ a hypothesis with error at most ϵ . Note that in reality m and t are not necessarily fixed, but in order for L to be a PAC-learning algorithm, they must be polynomial in all the relevant parameters. Therefore, the second source of error is taken care of.

1. (3pt) Explain the algorithm for generating a random bit above. Why does it produce a fair coin toss?

We assume the underlying distribution has some probability p of producing a positive label ($y_i = 1$) and $1 - p$ for negative ($y_i = 0$) where p is neither too close to 0 nor 1.

$$P(y_i \neq y_{i+1}) = p(1 - p) + (1 - p)p = 2p(1 - p)$$

Among the pairs $\{(0, 1), (1, 0)\}$ that appear, the chance of seeing $(0, 1)$ versus $(1, 0)$ is approximately 0.5. Thus, returning y_i is effectively returning 0 or 1 with equal probability, i.e., a fair coin toss. By discarding pairs with $y_i = y_{i+1}$, we ensure that we only use the pairs that differ. Although this might *slightly* reduce the total number of bits we can generate, each bit we do generate is closer to an unbiased coin toss.

2. (5pt) Find a lower bound on M such that if we choose h_0 or h_1 as the final hypothesis at the start, then with probability at least $1 - \delta/3$, we will have error at most ϵ .

The goal is to have a large enough sample size M so that if the deterministic procedure chooses h_0 or h_1 at the start, it will have true error at most ϵ with probability at least $1 - \delta/3$. Suppose the probability of a positive-labeled sample under distribution D is p_+ , then choosing h_0 has error $\text{Error}(h_0) = p_+$. We want $p_+ \leq \epsilon$ if we're forced to pick h_0 . If fewer than $\epsilon M/2$ of the M observed samples are positive, the algorithm picks h_0 . Let X_i be the indicator that sample i is labeled positive. Then the $\mathbb{E}[\sum_i X_i] = p_+ M$. If $p_+ > \epsilon$, the probability that we observe fewer than $\epsilon M/2$ positive samples is given by the Chernoff bound:

$$P\left(\sum_i X_i < \epsilon M/2\right) \leq e^{-M \cdot (\epsilon/2 - p_+)^2}$$

which is less than or equal to $\delta/3$ once M is large enough. By a similar argument, if $p_+ < 1 - \epsilon$, the probability that we see fewer than $\epsilon M/2$ negative examples is also small. By Chernoff bounds, we get a lower bound of

$$\frac{2}{\epsilon^2} \ln \frac{3}{\delta} \leq M$$

3. (7pt) Find a lower bound on M such that assuming we do not choose h_0 or h_1 as the final hypothesis at the start, with probability at least $1 - \delta/3$ we will generate at least t random bits using $M - m$ examples.

After deciding not to pick h_0 or h_1 , we want enough labeled examples left to produce the t random bits the learning algorithm needs. We want this to succeed with probability at least $1 - \delta/3$. Each pair of labels (y_i, y_{i+1}) is used to possibly generate 1 random bit if $y_i \neq y_{i+1}$. The chance that a pair differs is $2p(1 - p)$ if p is the probability of a positive label, call this p' . We need at least t successful pairs. This is similar to a binomial random variable with success probability p' in this case. The number of pairs that differ among N total pairs is therefore $p'N$. $P(\text{fewer than } t \text{ pairs differ}) \leq \delta/3$ if $\frac{t}{p'} \leq N$. If the algorithm needs m examples for training plus $2t$ labels for random bits, we require $M = m + \frac{2t}{p'}$ in order to ensure we do not run out of data for training. Hence, this gives a second bound on M , ensuring there are enough pairs to yield t random bits with high probability.