

Introduction to Machine Learning (67577)

Exercise 4 PAC & Ensemble Methods

Second Semester, 2021

Contents

1	Submission	2
2	PAC Learnability	2
3	VC-Dimension	2
4	Agnostic-PAC	3
5	Monotonicity	3
6	Theoretical Claim	3
7	Separate the Inseparable - Adaboost	4

1 Submission

- You should submit one file named **ex_4_FirstName_LastName.tar** containing the following files:
 - Ex4_Answers.pdf** - a pdf file which contains all your answers to both theoretical parts and practical parts.
 - adaboost.py**
 - ex4_tools.py**

2 PAC Learnability

- Let A be a learning algorithm, \mathcal{D} be any distribution, and our loss function is in the range $[0, 1]$ (e.g., the 0-1 loss). Prove that the following two statements are equivalent:
 - For every $\varepsilon, \delta > 0$, there exists $m(\varepsilon, \delta)$ such that $\forall m \geq m(\varepsilon, \delta)$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S)) \leq \varepsilon] \geq 1 - \delta$$

(b)

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathcal{A}(S))] = 0$$

Hint: Use Markov's inequality

- Sample Complexity of Concentric Circles in the Plane** Let $\mathcal{X} = \mathbb{R}^2$, $\mathcal{Y} = \{0, 1\}$ and let \mathcal{H} be the class of concentric circles in the plane, i.e., $\mathcal{H} = \{h_r : r \in \mathbb{R}_+\}$, where $h_r(x) = \mathbb{1}[\|x\|_2 \leq r]$. Prove that \mathcal{H} is PAC learnable and its sample complexity is bounded by

$$m_{\mathcal{H}}(\varepsilon, \delta) \leq \frac{\log(1/\delta)}{\varepsilon}$$

Note: Please do not use VC dimension arguments but instead prove the claim directly by showing a specific algorithm and analyzing its sample complexity.

Hint: Remember that for every ε ,

$$1 - \varepsilon \leq e^{-\varepsilon}$$

3 VC-Dimension

- Let $\mathcal{H} = \{h_1, \dots, h_N\}$ be finite class. show that:

$$VCdim(\mathcal{H}) \leq \lfloor \log_2 |\mathcal{H}| \rfloor$$

- Let $\mathcal{X} = \{0, 1\}^n$ and $\mathcal{Y} = \{0, 1\}$, for each $I \subseteq [n]$ define the parity function:

$$h_I(x) = \left(\sum_{i \in I} x_i \right) \bmod 2.$$

What is the VC-dimension of the class $\mathcal{H}_{\text{parity}} = \{h_I \mid I \subseteq [n]\}$? Prove your answer, you may use results from question 3.

- Given an integer k , let $([a_i, b_i])_{i=1}^k$ be any set of k intervals on \mathbb{R} and define their union $A = \cup_{i=1}^k [a_i, b_i]$. The hypothesis class $\mathcal{H}_{k\text{-intervals}}$ includes the functions:

$$h_A(x) = \begin{cases} 0 & x \notin A \\ 1 & x \in A \end{cases},$$

for all choices of k intervals. Find the VC-dimension of $\mathcal{H}_{k\text{-intervals}}$ and prove your answer. Show that if we let A be any finite union of intervals (i.e. k is unlimited), then the resulting class $\mathcal{H}_{\text{intervals}}$ has VC-dimension ∞ .

6. **Boolean Conjunctions** Let $\mathcal{X} = \{0, 1\}^d$ and $\mathcal{Y} = \{0, 1\}$, and assume $d \geq 2$. Each sample $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ consists of an assignment to d boolean variables (\mathbf{x}) and a label (y). For each boolean variable $x_k, k \in [d]$, there are two literals: x_k and $\bar{x}_k = 1 - x_k$. The class \mathcal{H}_{con} is defined by boolean conjunctions over any subset of these $2d$ literals. Compute the VC dimension of \mathcal{H}_{con} and prove your answer.

4 Agnostic-PAC

7. Prove that if \mathcal{H} has the uniform convergence property with function $m_{\mathcal{H}}^{UC} : (0, 1)^2 \rightarrow \mathbb{N}$, then \mathcal{H} is Agnostic-PAC learnable with sample complexity $m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}(\epsilon/2, \delta)$.
8. Let \mathcal{H} be a hypothesis class over a domain $Z = \mathcal{X} \times \{\pm 1\}$, and consider the 0-1 loss function. Assume that there exists a function $m_{\mathcal{H}}$, for which it holds that for every distribution \mathcal{D} over Z there is an algorithm \mathcal{A} with the following property: when running \mathcal{A} on $m \geq m_{\mathcal{H}}$ i.i.d. examples drawn from \mathcal{D} , it is guaranteed to return, with probability at least $1 - \delta$, a hypothesis $h_S : \mathcal{X} \rightarrow \{\pm 1\}$ with $L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$. Is \mathcal{H} agnostic PAC learnable? Prove or show a counter example.

5 Monotonicity

9. **Of Sample Complexity** Let \mathcal{H} be a hypothesis class for a binary classification task. Suppose that \mathcal{H} is PAC learnable and its sample complexity is given by $m_{\mathcal{H}}(\cdot, \cdot)$. Show that $m_{\mathcal{H}}$ is monotonically non-increasing in each of its parameters. That is, show that given $\delta \in (0, 1)$, and given $0 < \epsilon_1 \leq \epsilon_2 < 1$, we have that $m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta)$. Similarly, show that given $\epsilon \in (0, 1)$, and given $0 < \delta_1 \leq \delta_2 < 1$, we have that $m_{\mathcal{H}}(\epsilon, \delta_1) \geq m_{\mathcal{H}}(\epsilon, \delta_2)$.
10. **of VC-Dimension** Let \mathcal{H}_1 and \mathcal{H}_2 be two classes for binary classification, such that $\mathcal{H}_1 \subseteq \mathcal{H}_2$. Show that $VC - \dim(\mathcal{H}_1) \leq VC - \dim(\mathcal{H}_2)$.

6 Theoretical Claim

11. **(optional)** Let \mathcal{X} be a sample space and $\mathcal{Y} = \{\pm 1\}$. Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a hypothesis class. For $C \subseteq \mathcal{X}$, recall the notation \mathcal{H}_C for the restriction of \mathcal{H} to the subset C . Define the function $\tau_m(\mathcal{H}) : \mathbb{N} \rightarrow \mathbb{N}$ corresponding to \mathcal{H} to be

$$\tau_{\mathcal{H}}(m) := \max \left\{ |\mathcal{H}_C| \mid C \subseteq \mathcal{X}, |C| = m \right\}.$$

- (a) Explain, in your own words, the meaning of $\tau_{\mathcal{H}}$.
- (b) Suppose that $VCdim(\mathcal{H}) = \infty$. Find an expression for the value of $\tau_{\mathcal{H}}(m)$ for $m \in \mathbb{N}$.
- (c) Now suppose that $VCdim(\mathcal{H}) = d$. Find an expression for the value of $\tau_{\mathcal{H}}(m)$ for $m \leq d$.
- (d) You will now prove the following important result: suppose that $VCdim(\mathcal{H}) = d$ and let $m > d$. Then

$$\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d} \right)^d,$$

where e is the natural logarithm base. You'll do this in three steps:

- i. Using induction, show that for any finite $C \subseteq \mathcal{X}$,

$$|\mathcal{H}_C| \leq \left| \{B \subseteq C \mid \mathcal{H} \text{ shatters } B\} \right|.$$

Hint: in the induction step divide \mathcal{H}_C to two groups. one of them can be $\mathcal{H}_{C'}$ when $C' = \{c_2, \dots, c_m\}$.

- ii. Explain in your own words the meaning of this inequality.
 iii. Show that, for any finite $C \subset \mathcal{X}$, we have

$$\left| \{B \subset C \mid \mathcal{H} \text{ shatters } B\} \right| \leq \sum_{k=0}^d \binom{m}{k}$$

- iv. Use the following inequality (which you are not required to prove)

$$\sum_{k=0}^d \binom{m}{k} \leq \left(\frac{em}{d}\right)^d$$

to finish the proof that $\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$.

- (e) If $m = d$, does the inequality $\tau_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$ hold? If it does hold, is it tight?
 (f) Characterize in words the behavior of $\tau_{\mathcal{H}}(m)$ for $m \leq VCdim(\mathcal{H})$ and for $m > VCdim(\mathcal{H})$. Can you use your characterization to offer an alternative definition of the VC-dimension $VCdim(\mathcal{H})$?

7 Separate the Inseparable - Adaboost

- As we saw in class, one of the main motivations for boosting is computational complexity. Try to solve the following questions in a reasonable runtime. Two tips to help reducing runtime:
 - Try to **Avoid for loops**. If you are not familiar with the concept of vectorization in numpy, this is the time for you to read about it. You can find many tutorials on the web, for example [this](#). (loops are allowed as long as your code runs in reasonable time).
 - You can use [line_profiler](#). It is a tool that helps in finding the slower parts of your code. After you located the bottleneck try to think if you can accelerate it.
- The file `adaboost.py` contains a template for implementation of the adaboost classifier. Fill in the template and implement the Adaboost algorithm as learned in class, where you assume `WL` is a weak-learner class that can be used as follows:
 - `h = WL(D,X,y)` - constructs a weak classifier trained on X, y weighted by D .
 - `h.predict(X)` - returns the classifier's prediction on a set X .
 In `ex4_tools` you are provided with such an implementation for a Decision Stump classifier, called `DecisionStump`. You may add methods as you see fit, but please implement the functions that are declared in the template.
- In `ex4_tools` you are provided with the function `generate_data`. Use it to generate 5000 samples without noise (i.e. `noise_ratio=0`). Train an Adaboost classifier over this data. Use the `DecisionStump` weak learner mentioned above, and $T = 500$. Generate another 200 samples without noise ("test set") and plot the training error and test error, as a function of T . Plot the two curves on the same figure.

14. Plot the decisions of the learned classifiers with $T \in \{5, 10, 50, 100, 200, 500\}$ together with the test data. You can use the function `decision_boundaries` together with `plt.subplot` for this purpose.
15. Out of the different values you used for T , find \hat{T} , the one that minimizes the test error. What is \hat{T} and what is its test error? Plot the decision boundaries of this classifier together with the training data.
16. Look into the AdaBoost: Take the weights of the samples in the last iteration of the training (D^T). Plot the training set with size proportional to its weight in D^T , and color that indicates its label (again, you can use `decision_boundaries`). Oh! we cannot see any point! the weights are too small... so we will normalize them: $D = D / \text{np.max}(D) * 10$. What do we see now? can you explain it?
17. Repeat 13,14,15,16 with noised data. Try `noise_ratio=0.01` and `noise_ratio=0.4`.
 - Add all the graphs to the pdf.
 - Describe the changes.
 - Explain 13 in terms of the bias complexity tradeoff.
 - Explain the differences in 15.