

Homework 2: Oct 28th, 2018

Due: Nov 11th, 2018 (See the submission guidelines in the course web site)

Theory Questions

1. **(12 points) Reject Option in Classifiers.** In this question we will consider a classification setup with an additional reject option. This can be useful in cases of uncertainty, and when the cost for rejection is less than the cost of misclassifying the input, it may be the optimal action. Given a data point $\langle x, y \rangle$, the output \hat{y} of a classifier is defined over the set $\{1, \dots, L, L + 1\}$, where L is the number of classes and $L + 1$ is the reject option. Let $\lambda \in \mathbb{R}$, the loss function is defined as follows:

$$\Delta(y, \hat{y}) = \begin{cases} 0 & \text{if } \hat{y} = y \\ \lambda & \text{if } \hat{y} = L + 1 \\ 1 & \text{otherwise} \end{cases}$$

In other words, the loss is 0 for correct classification, λ for rejection and 1 for misclassification.

- (a) **(8 points)** Show that the minimum expected loss is obtained by a classifier that returns the class $i \in \{1, \dots, L\}$ if i is the most probable class and $\mathbb{P}[Y = i | X = x] \geq 1 - \lambda$, and rejects otherwise.
- (b) **(4 points)** Describe what happens as λ is increased from 0 to 1.
2. **(16 points) Singletons** (Section 3.5, Ex.2 in the course textbook). Let \mathcal{X} be a discrete domain, and let $\mathcal{H}_{\text{Singleton}} = \{h_z : z \in \mathcal{X}\} \cup \{h^-\}$, where for each $z \in \mathcal{X}$, h_z is the function defined by $h_z(x) = 1$ if $x = z$ and $h_z(x) = 0$ if $x \neq z$. h^- is simply the all-negative hypothesis, namely, $\forall x \in \mathcal{X}, h^-(x) = 0$.
- The realizability assumption here implies that the true hypothesis f labels negatively all examples in the domain, perhaps except one.
- (a) **(8 points)** Describe an algorithm that implements the ERM rule for learning $\mathcal{H}_{\text{Singleton}}$ in the realizable setup.
- (b) **(8 points + 5 bonus points)** Show that $\mathcal{H}_{\text{Singleton}}$ is PAC learnable, and provide $N(\epsilon, \delta)$. A direct proof for the algorithm, without using VC dimension, will receive up to 5 additional bonus points.
3. **(10 points) PAC in Expectation.** Consider learning in the realizable case. We say an hypothesis class \mathcal{H} is *PAC learnable in expectation* if there exists a learning

algorithm A and a function $N(a) : (0, 1) \rightarrow \mathbb{N}$ such that $\forall a \in (0, 1)$ and for any distribution P , given a sample set S , such that $|S| \geq N(a)$ it holds that,

$$\mathbb{E}[e_P(A(S))] \leq a$$

Show that \mathcal{H} is PAC learnable *if and only if* \mathcal{H} is PAC learnable in expectation (Hint: Use Markov's inequality and refer to derivations between equations 3.13-3.14 in the lecture scribes about VC).

4. **(10 points) Union of Intervals.** Determine the VC-dimension of the subsets of the real line formed by the union of k intervals (see question 1 of the programming assignment for a formal definition of \mathcal{H}).
5. **(10 points) Structural Risk Minimization.** Let \mathcal{H} be a countable hypothesis class, that is, \mathcal{H} can be written as $\mathcal{H} = \bigcup_{i \in \mathbb{N}} \{h_i\}$. Let $w : \mathcal{H} \rightarrow [0, 1]$ be a function such that $\sum_{h \in \mathcal{H}} w(h) \leq 1$. We refer to w as a *weight function* over the hypotheses which reflects the prior for each hypothesis.

Show that with probability $1 - \delta$ over the choice $S \sim P$ ($|S| = n$)

$$\forall h \in \mathcal{H}, e_P(h) \leq e_S(h) + \sqrt{\frac{1}{2n} \ln \frac{2}{\delta \cdot w(h)}}$$

Hint: use the uniform convergence property for each hypothesis class $\{h_i\}$ of size 1.

Programming Assignment

1. **(42 points) Union of Intervals.** In this question, we will study the hypothesis class of a finite union of disjoint intervals, and the properties of the ERM algorithm for this class.

To review, let the sample space be $\mathcal{X} = [0, 1]$ and assume we study a binary classification problem, i.e. $\mathcal{Y} = \{0, 1\}$. We will try to learn using an hypothesis class that consists of k intervals. More explicitly, let $I = \{[l_1, u_1], \dots, [l_k, u_k]\}$ be k disjoint intervals, such that $0 \leq l_1 \leq u_1 \leq l_2 \leq u_2 \leq \dots \leq u_k \leq 1$. For each such k disjoint intervals, define the corresponding hypothesis as

$$h_I(x) = \begin{cases} 1 & \text{if } x \in [l_1, u_1], \dots, [l_k, u_k] \\ 0 & \text{otherwise} \end{cases}$$

Finally, define \mathcal{H}_k as the hypothesis class that consists of all hypotheses that correspond to k disjoint intervals:

$$\mathcal{H}_k = \{h_I | I = \{[l_1, u_1], \dots, [l_k, u_k]\}, 0 \leq l_1 \leq u_1 \leq l_2 \leq u_2 \leq \dots \leq u_k \leq 1\}$$

We note that $\mathcal{H}_k \subseteq \mathcal{H}_{k+1}$, since we can always take one of the intervals to be of length 0 by setting its endpoints to be equal. We are given a sample of size $m = \langle x_1, y_1 \rangle, \dots, \langle x_n, y_m \rangle$. Assume that the points are sorted, so that $0 \leq x_1 < x_2 < \dots < x_m \leq 1$.

Submission guidelines

- Download the files `skeleton.py` and `intervals.py` from Moodle. You should implement only the missing code in `skeleton.py`, as specified in the following questions. In every method description, you will find specific details on its input and return values.
- Your code should be written with python 3.
- Make sure to comment out / remove any code which halts the code execution, such as matplotlib popup.
- Your submission should include exactly two files: `assignment2.py`, `intervals.py`.

Explanation on intervals.py

The file `intervals.py` includes a function that implements an ERM algorithm for \mathcal{H}_k . Given a sorted list $\mathbf{xs}=[x_1, \dots, x_m]$, the respective labeling $\mathbf{ys}=[y_1, \dots, y_m]$ and k , the given function `find_best_interval` returns a list of up to k intervals and their error count on the given sample. These intervals have the smallest empirical error count possible from all choices of k intervals or less.

Note that in sections (d)-(f) you will need to use this function for large values of m . Execution in these cases could take time (more than 10 minutes for experiment), so plan ahead.

- (a) **(7 points)** Assume that the true distribution $P[x, y] = P[y|x] \cdot P[x]$ is: x is distributed uniformly on the interval $[0, 1]$, and

$$P[y = 1|x] = \begin{cases} 0.8 & \text{if } x \in [0, 0.2] \cup [0.4, 0.6] \cup [0.8, 1] \\ 0.1 & \text{if } x \in [0.2, 0.4] \cup [0.6, 0.8] \end{cases}$$

and $P[y = 0|x] = 1 - P[y = 1|x]$.

Implement the method `sample_from_D` that draws m pairs of (x, y) according to the distribution P . Then, implement the method `draw_sample_intervals`, that draws a sample of size $m = 100$ using `sample_from_D` and creates a plot:

- i. Plot the points and their label (have the y axis in range $-0.1, 1.1$ for clarity of presentation).
 - ii. Mark the lines $x = 0.2, 0.4, 0.6, 0.8$ clearly on the plot.
 - iii. Run the `find_best_interval` function on your sample with $k = 3$, and plot the intervals clearly.
- (b) **(7 points)** Note that here, we know the true distribution P , so for every given hypothesis $h \in \mathcal{H}_k$, we can calculate $error(h)$ precisely. What is the hypothesis with the smallest error?
- (c) **(7 points)** Implement the method `experiment_m_range_erm` that given a list of intervals, calculates the error of the respective hypothesis. For $k = 3$, $m = 10, 15, 20, \dots, 100$, perform the following experiment $T = 100$ times: (i) Draw a sample of size m and run the ERM algorithm on it; (ii) Calculate the empirical error for the returned hypothesis; (iii) Calculate the true error for the returned hypothesis. Plot the average empirical and true errors, averaged across the T runs, as a function of m . Discuss the results. Do the empirical and true error decrease or increase in m ? Why?
- (d) **(7 points)** Implement the method `experiment_k_range_erm`, to perform the following experiment. Draw a data set of $m = 1500$ samples. Find the best ERM hypothesis for $k = 1, 2, \dots, 10$, and plot the empirical and true errors as a function of k . How does the error behave? Define k^* to be the k with the smallest empirical error for ERM? Does this mean the hypothesis with k^* intervals is a good choice?
- (e) **(7 points)** Now we will use the principle of structural risk minimization (SRM), to search for a k that gives good test error. Let $\delta = 0.1$:
- Use your results from question 4 in the theoretical part and the VC confidence bound to construct a penalty function.
 - Implement the method `experiment_k_range_srm`, to perform the following experiment. Draw a data set of $m = 1500$ samples, run the experiment in (d) again, but now plot two additional lines as a function of k : 1) the penalty for the best ERM hypothesis and 2) the sum of penalty and empirical error.
 - What is the best value for k in each case? is it better than the one you chose in (d)?

- (f) **(7 points)** Here we will use holdout-validation to search for a k that gives good test error. Implement the method `cross_validation` to perform the following. Draw a data set of $m = 1500$ samples and use 20% for a holdout-validation. Choose the best hypothesis based on 3 experiments. Discuss how close this gets you to finding the hypothesis with optimal true error.