

Seminar 5

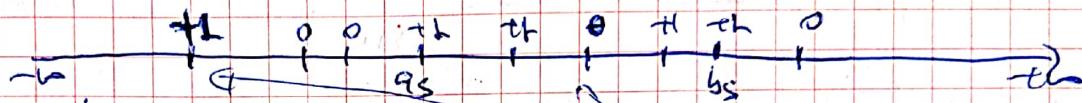
① $\mathcal{H}_{\text{intervals}} = \mathcal{H}_{\text{vec}} = \{ h_{a,b} : \mathbb{R} \rightarrow \mathbb{R}, h_{a,b} = \mathbb{1}_{[a,b]}, h_{a,b}(x) = \begin{cases} 1, & x \in [a,b] \\ 0, & \text{otherwise}, a, b \in \mathbb{R} \end{cases}\}$

Consider a training set S of size m .

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m) | x_i \in \mathbb{R}, y_i \in \{0, 1\}, 1 \leq i \leq m\}$$

Propose an implementation of the ERM_{ag} learning rule in the agnostic case that runs in $O(m^2)$ to find a hypothesis h_{a_*, b_*} with the smallest empirical risk.

Example:



The solution is h_{a_*, b_*} with loss $\frac{2}{9}$. (There are 2 misclassified points out of 9.)

Observations

1. We are in the agnostic case:
 - it might be the case that there is no labelling function but instead we are dealing with a distribution (some point might have different labels)
 - if there is a labelling function it might not be in $\mathcal{H}_{\text{intervals}}$
2. If all points are negative we should return an interval not containing any point in S .
3. If all points are positive should return an interval containing all points in S .

We will first sort the training set S in ascending order of x 's. So we

obtain $S = \{(x_{(1)}, y_{(1)}), (x_{(2)}, y_{(2)}), \dots, (x_{(m)}, y_{(m)})\}$

$$\text{with } x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)}$$

As we are in the agnostic case we can have $x_{(1)} = x_{(i+1)}$ and

$$y_{(i)} \neq y_{(i+1)}$$

Consider the set \mathcal{Z} containing the values of x with no repetition

$$\mathcal{Z} = \{z_1, z_2, \dots, z_n\}$$

$$z_{(1)} x_{(1)} < z_2 < \dots < z_{(m)} x_{(m)}, n \leq m.$$

If all initial x values are different than $z_{(1)} x_{(1)}, \dots, z_{(m)} x_{(m)}$, $n = m$.

Code of the implementation of ERM_{ag}

1. If all $y_i = 0$ return an interval not containing any point $X: [z_{(1)}, z_{(m)}]$.

2. Consider all possible intervals $Z_{ij} = [z_i, z_j]$ $i = \overline{1, m}, j = \overline{i, m}$

There are $m(m+1) + (m-1) + \dots + 1 = \frac{m(m+1)}{2}$ such intervals.

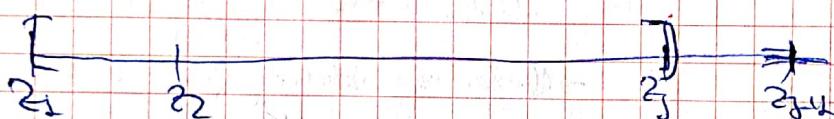
Consider to determine the interval $Z^* = Z_{i^*, j^*}$ with the smallest empirical risk Z_{i^*, j^*} $\approx \arg\min_{i=1, j=i} \text{Loss}(Z_{ij})$

How to compute very fast $\text{Loss}(Z_{ij})$? Use dynamic programming!

$$\text{loss}(Z_{ij}) = \frac{\# \text{negative points inside } Z_{ij}}{m} + \frac{\# \text{positive points outside } Z_{ij}}{m}$$

Key observation: $\text{loss}(Z_{i+1, j})$ can be easily computed based on $\text{loss}(Z_{ij})$

Simple case: there is just one point (x_k, y_k) in the training set S such that $x_k \in Z_{i+1, j}$



If $y_k = +1$ then $\text{loss}(Z_{i, j+1}) = \text{loss}(Z_{ij}) - \frac{1}{m}$ (the loss decreases)

If $y_k = 0$ then $\text{loss}(Z_{i, j+1}) \geq \text{loss}(Z_{ij}) + \frac{1}{m}$. (the loss increases)

General case (in the agnostic scenario)

We have multiple points $x_{k_1}, x_{k_2}, \dots, x_{k_l} \in Z_{i, j+1}$ (l points)

Then some of the points will have label $+ = p_{j+1}$

Some of the points will have label $0 = n_{j+1} - p_{j+1}$

In this case $\text{loss}(Z_{i, j+1}) = \text{loss}(Z_{ij}) - \frac{p_{j+1}}{m} + \frac{n_{j+1}}{m}$

because p_{j+1} points will be labeled correctly now

n_{j+1} points will be labeled incorrectly now

(if $l = 1$ we have $p_{j+1} + n_{j+1} = 1$ we have just one point labeled positive or negative)

Efficient implementation of the ERM_{0/1} rule for Blankwals

1. Sort S and obtain $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)}$. Build the set $Z = \{x_1, x_2, \dots, x_m\}$ containing value x without repetition.

$$Z = \{x_1, x_2, \dots, x_m\} \quad x_{(1)} < x_2 < \dots < x_m = x_{(m)}$$

2. Check if all $y_i \in \{-1, 1\}$ have defined. If so return $\{a_S, b_S\}$, where $a_S = z_{(2)}, b_S = z_{(1)}$. Compute $P = \sum_{i=1}^m y_i$ (the positive examples)

3. For $j = \overline{1, m}$

compute values $p_j = \# \text{points } x_i = z_j \text{ with label } y_i = 1$

$n_j = \# \text{points } x_i = z_j \text{ with label } y_i = 0$

$$\leftarrow \min_error = \frac{m}{m} = 1, i^* = [], j^* = []$$

for $i = \overline{1, m}$

for $j = \overline{1, m}$

$$z_{ij} = [z_i, p_j]$$

$$y_j (y_j = -1)$$

$$\text{loss}(z_{ij}) = \frac{P - p_j + n_j}{m}$$

if $\text{loss}(z_{ij}) < \min_error$

else

$$\text{loss}(z_{ij}) = \text{loss}(z_{i, j-1}) + \frac{n_j - p_j}{m}$$

if $\text{loss}(z_{ij}) < \min_error$

$$\min_error = \text{loss}(z_{ij})$$

$$i^* = i$$

$$j^* = j$$

5. Return i^*, j^*

Complexity: 1 sorting $O(m \log m)$

2 computing $P = O(m)$

3 computing $p_j, n_j = O(m)$

4 computing $\text{loss}(z_{ij})$ = constant time $O(1)$

② If $\mathcal{H} = \{h_{a,b,s}: \mathbb{R} \rightarrow \{-1, +1\}\}$, $a, b, s \in \{-1, +1\}$, where

$$h_{a,b,s}(x) = \begin{cases} 1, & x \in [a, b] \\ -1, & x \notin [a, b] \end{cases}$$

a) reduct case

There exists a function $h_{a^*, b^*, s^*} \in \mathcal{H}$ that matches the training points

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

$$y_i = h_{a^*, b^*, s^*}(x_i)$$

We can have the following possibilities for examples appearing in S

t t t t t (only positive examples)

- - - - - (only negative examples)

t t t - - t t t

- - - t t t - - -

t t t t - - -

- - - - - t t t t

Consider the following algorithm

Initialization: $a_+ = -\infty, a_- = \infty$

$$b_+ = +\infty, b_- = -\infty$$

Compute $a_+ = \min_{\substack{i=1, m \\ y_i=+1}} x_i$, if there is two x_i with $y_i=+1$ then $a_+ = -\infty$

$$b_+ = \max_{\substack{i=1, m \\ y_i=+1}} x_i, \quad | | - | | -$$

$$b_+ = +\infty$$

$$a_- = \min_{\substack{i=1, m \\ y_i=-1}} x_i, \quad - ? | | -$$

$$a_- = -\infty$$

$$b_- = \max_{\substack{i=1, m \\ y_i=-1}} x_i, \quad - | | -$$

$$b_- = -\infty$$

If $a_+ < a_-$ return $h_{a_-, b_-, -1}$
 Else return h_{a^*, b^*, s^*}

b) agnostic case

Can think off \mathcal{H} as signed intervals $\cup \mathcal{H}^+$ intervals

$$\mathcal{H}^+ \text{ intervals} = \{ h_{a,b} : R \rightarrow \{-1, +1\}, a \leq b \}$$
$$h_{a,b}(x) = \begin{cases} 1 & x \in [a, b] \\ -1 & x \notin [a, b] \end{cases}$$

$$\mathcal{H}^- \text{ intervals} = \{ h_{a,b} : R \rightarrow \{-1, +1\}, a \leq b \}$$
$$h_{a,b}(x) = \begin{cases} -1 & x \in [a, b] \\ 1 & x \notin [a, b] \end{cases}$$

Use the algorithm in exercise 1 (efficient implementation of the ERM rule) and run it for \mathcal{H}^+ intervals as \mathcal{H} intervals.

(Observe the hypothesis $h_{a,b}^{+}$ and $h_{a,b}^{-}$)

Choose the one with the minimum empirical error.

③ A algorithm with the following property:

$\exists \delta_0 \in (0, 1)$ and $m_{\mathcal{H}} : (0, 1) \rightarrow N$ s.t. for every $\varepsilon \in (0, 1)$ if $m \geq m_{\mathcal{H}}(\varepsilon)$ then for every D distribution it holds

$$\Pr_{S \sim D^m} (L_D(A(S)) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon) \geq 1 - \delta_0$$

Suggest a procedure based on algorithm A that learns \mathcal{H} in the agnostic PAC setting and has a sample complexity of

$$m_{\mathcal{H}}(\varepsilon, \delta) \leq k \cdot m_{\mathcal{H}}\left(\frac{\varepsilon}{2}\right) + \lceil \frac{2 \cdot \log(4k/\delta)}{\varepsilon^2} \rceil,$$

$$k = \lceil \log \frac{\delta}{2} / \log \delta_0 \rceil$$

Def of Agnostic PAC: \mathcal{H} is agnostic PAC if there exists a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow N$

and a learning algorithm A' with the following property: $\forall \varepsilon > 0, \forall \delta > 0, \forall D$ distribution over $Z = \{0, 1\}$ when we run the algorithm A' on a learning set S of $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ examples sampled i.i.d. from D it returns $h_S = A'(S)$ s.t.

$$\Pr_{S \sim D^m} (L_D(h_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \varepsilon) \geq 1 - \delta$$

This is equivalent to:

$$P_{S \sim D^{\text{true}}} (L_D(h_S) > \min_{h \in \mathcal{H}} L_D(h) + \frac{\epsilon}{2}) < \delta$$

$\underbrace{\min_{h \in \mathcal{H}} L_D(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2}}$

Follow the indications.

Let $\epsilon, \delta \in (0, 1)$. Pick k 'samples' S_1, S_2, \dots, S_k of size $m_k(\frac{\epsilon}{2})$

Use the property of the algorithm & given.

$$\forall i=1, k \quad A(S_i) \approx h_i$$

$$P_{S_i \sim D^{\text{true}}(\frac{\epsilon}{2})} (L_D(h_i) \leq \min_{h \in \mathcal{H}} L_D(h) + \frac{\epsilon}{2}) \geq 1 - \delta$$

$$\Rightarrow P_{S \sim D^{\text{true}}(\frac{\epsilon}{2})} (L_D(h_i) > \min_{h \in \mathcal{H}} L_D(h) + \frac{\epsilon}{2}) < \delta.$$

(the probability of having a bad h_i).

The probability that all h_1, \dots, h_k are bad is given by,

$$P(L_D(h_1) > \min_{h \in \mathcal{H}} L_D(h) + \epsilon_1 \text{ and } L_D(h_2) > \min_{h \in \mathcal{H}} L_D(h) + \epsilon_2 \text{ and } \dots) \leq \delta^k$$

~~Take~~ $\delta^k \leq \delta$ Find k such that $\delta^k \leq \delta$

$$\Leftrightarrow k \cdot \log \delta \leq \log \frac{\delta}{2} \quad | : \log \delta$$

$$k \geq \lceil \frac{\log \delta - \log 2}{\log \delta} \rceil$$

Consider $H'_2 \{h_1, h_2, \dots, h_k\}$

H' finite, apply Corollary (4.6)

if $m \geq m_{H'}(\frac{\epsilon}{2}, \frac{\delta}{2}) \leq \lceil \frac{2 \log(4K/\delta)}{\epsilon^2} \rceil$ we have that

$$P_{S \sim D^{\text{true}}(\frac{\epsilon}{2}, \frac{\delta}{2})} (L_D(h_{k+1}) > \min_{h \in H} L_D(h) + \frac{\epsilon}{2}) < \frac{\delta}{2}$$

See that $\lceil \frac{2 \log(4K/\delta)}{\epsilon^2} \rceil$ example

So: $L_D(h_{K+1}) > \min_{h \in H} L_D(h) + \frac{\epsilon}{2}$ if either we have

A: all h_i are bad $\Rightarrow L_D(h_i) > \min_{h \in H} L_D(h) + \frac{\epsilon}{2}$

B: h_{K+1} is bad $\Leftrightarrow L_D(h_{K+1}) > \min_{h \in H} L_D(h) + \frac{\epsilon}{2}$

$$P(A \cup B) \leq P(A) \cup P(B) \leq \frac{d}{\epsilon} + \frac{d}{\epsilon} = d.$$

So, Place $m = K \cdot m_0(\frac{\epsilon}{2}) + \lceil \frac{-\log(4K/\delta)}{\epsilon_2} \rceil$, $K \geq \lceil \frac{\ln \delta - \ln 2}{\ln \delta_0} \rceil$

$P(L_D(h_{K+1}) > \min_{h \in H} L_D(h) + \frac{\epsilon}{2}) < \delta$ v.

$(S_1, S_2, \dots, S_K, S_{K+1}) \sim D^m$
 h_1, h_2, \dots, h_K bad