

IML Ex4

Eliyahu Strugo

June 5, 2020

PAC Learnability

1.Solution:

- (a) \Rightarrow (b):

Assume that for any $\epsilon, \delta > 0$, there exists $m(\epsilon, \delta)$ s.t $\forall m \geq m(\epsilon, \delta)$ it holds that:

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

Let $\epsilon > 0$ and $\delta = \frac{\epsilon}{2}$.

$$\alpha := \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) | L_{\mathcal{D}}(h_S) > \frac{\epsilon}{2}] \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \frac{\epsilon}{2}]$$

$$\beta := \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) | L_{\mathcal{D}}(h_S) \leq \frac{\epsilon}{2}] \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) \leq \frac{\epsilon}{2}]$$

By the Law of Total Expectation:

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] = \alpha + \beta$$

$$- \alpha : \text{Since } L_{\mathcal{D}}(h_S) \leq 1 \text{ then } \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) | L_{\mathcal{D}}(h_S) > \frac{\epsilon}{2}] \leq 1$$

$$\text{and there exists } m(\frac{\epsilon}{2}, \frac{\epsilon}{2}) \text{ s.t } \forall m \geq m(\frac{\epsilon}{2}, \frac{\epsilon}{2}) \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \frac{\epsilon}{2}] \leq \frac{\epsilon}{2} \Rightarrow \alpha \leq \frac{\epsilon}{2}$$

$$- \beta : \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) | L_{\mathcal{D}}(h_S) \leq \frac{\epsilon}{2}] \leq \frac{\epsilon}{2} \text{ and } \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) \leq \frac{\epsilon}{2}] \leq 1 \Rightarrow \beta \leq \frac{\epsilon}{2}$$

Hence, $\forall m \geq m(\frac{\epsilon}{2}, \frac{\epsilon}{2})$ it holds that $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] = \alpha + \beta \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon$ and therefore:

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

- (b) \Rightarrow (a):

Assume that $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] \xrightarrow{m \rightarrow \infty} 0 \Rightarrow \forall \epsilon > 0 \exists m(\epsilon) \in \mathbb{N} \text{ s.t } \forall m \geq m(\epsilon) :$

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)] \stackrel{*}{=} |\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]| < \epsilon$$

* Since $0 \leq L_{\mathcal{D}}(h_S)$ then $0 \leq \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]$.

Let $\epsilon > 0$. By the Markov's inequality it holds that:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \epsilon] \leq \frac{\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]}{\epsilon}$$

Now ,Let $\delta > 0$ then $\epsilon\delta > 0$ and $\exists m(\epsilon, \delta) \in \mathbb{N} \text{ s.t } \forall m \geq m(\epsilon, \delta) :$

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \epsilon] \leq \frac{\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S)]}{\epsilon} < \frac{\epsilon\delta}{\epsilon} = \delta$$

Hence:

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

2.Solution:

Given training set $S = ((x_1, y_1), \dots, (x_m, y_m))$, define $A = \{r \in \mathbb{R}_+ | \forall (x, 1) \in S \text{ it holds that } \|x\|_2 \leq r\}$.

Our algorithm is ERM which finds a radius r_{alg} s.t:

$$r_{alg} = \begin{cases} \min\{A\} & A \neq \emptyset \\ 0 & \text{else} \end{cases}$$

and returns a hypothesis $h_S = A(S)$. Fix a distribution \mathcal{D} over \mathcal{X} . Since we assuming realizability then \mathcal{H} contains the true hypothesis. Fix true hypothesis $h_{r^*} \in \mathcal{H} \Rightarrow L_{\mathcal{D}}(h_{r^*}) = 0$. By the definition of r_{alg} it holds that:

$$r_{alg} < r^*$$

Let $x \in \mathcal{X}$. if $\|x\| \leq r_{alg}$ or $r^* < \|x\|$ then the algorithm is correct and if $r_{alg} < \|x\| \leq r^*$ the algorithm will be incorrect. Therefore, we want to examine the case where the algorithm is incorrect.

Let $\epsilon > 0$.

If $\mathcal{D}(\{x : \|x\| \leq r^*\}) < \epsilon \Rightarrow \mathcal{D}(\{x : r_{alg} < \|x\| \leq r^*\}) < \epsilon \Rightarrow \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) \leq \epsilon] = 1$ and we are done, otherwise $\mathcal{D}(\{x : \|x\| \leq r^*\}) \geq \epsilon$.

Let us define r' s.t $\mathcal{D}(\{x : r' \leq \|x\| \leq r^*\}) \leq \epsilon$ and define $B = \{x : r' \leq \|x\| \leq r^*\}$. Since $A(S)$ picks the minimal radius that contains all positive instances then:

$$S \cap B \neq \emptyset \Rightarrow r' \leq r_{alg}$$

Which means:

$$\{x : r_{alg} < \|x\| \leq r^*\} \subseteq B$$

$$\Rightarrow \mathcal{D}(\{x : r_{alg} < \|x\| \leq r^*\}) \leq \mathcal{D}(B) \leq \epsilon$$

And we get $S \cap B \neq \emptyset \Rightarrow L_D(h_S) \leq \epsilon$, which holds iff $L_D(h_S) > \epsilon \Rightarrow S \cap B = \emptyset$, and then:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \epsilon] \leq \mathbb{P}_{S \sim \mathcal{D}^m} [S \cap B = \emptyset] = \mathbb{P}_{S \sim \mathcal{D}^m} [\forall (x, y) \in S, (x, y) \notin B] \stackrel{i.i.d}{=} \prod_{i=1}^m \mathbb{P}_{x_i, y_i \sim \mathcal{D}^m} [(x_i, y_i) \notin B]$$

Moreover, $\mathcal{D}(B) \leq \epsilon \Rightarrow \mathcal{D}(B \cap S) \leq \epsilon \Rightarrow \forall i \in [m], \mathbb{P}_{x_i, y_i \sim \mathcal{D}^m} [\{x_i, y_i\} \in B] \leq \epsilon$ then:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(h_S) > \epsilon] \leq \prod_{i=1}^m \mathbb{P}_{x_i, y_i \sim \mathcal{D}^m} [(x_i, y_i) \notin B] \leq (1 - \epsilon)^m \leq e^{-m\epsilon}$$

And $e^{-m\epsilon} \leq \delta \iff \frac{1}{\epsilon} \ln\left(\frac{1}{\delta}\right) \leq m$.

VC dimension

3.Solution:

There is shattered sample of size d :

Let $C = \{v_1, \dots, v_d\}$ s.t $v_1, \dots, v_d \in \mathbb{R}^d$ and $\forall i, j \in [d], v_i(j) = 1 - \delta_{i,j}$, where $v_i^{(j)}$ is the j -th coordinate of v_i and $\delta_{i,j}$ is Kronecker delta.

Let $y = (y_1, \dots, y_d) \in \{0, 1\}^d$ and define $h_y = \bigwedge_{i: y_i=0} x_i$ if $y \neq 1_{\mathbb{R}^d}$ else $h_y = 1$. For any $i \in [d]$ it holds that $h_y(v_i) = y_i$. Thus, C is shattered $\Rightarrow VCdim(\mathcal{H}) \geq d$.

There is no shattered sample of size $d+1$:

Assume towards contradiction that there exists $C = \{v_1, \dots, v_{d+1}\}$ that is shattered by \mathcal{H} , which means that points in C can have all possible labeling $\Rightarrow \forall i, j \in [d+1]$ there exists $h_i \in \mathcal{H}$ s.t $h_i(v_j) = 1 - \delta_{i,j}$.

If $h_i(v_j) = 1 - \delta_{i,j}$ then the conjunction h_i must contain some literal $l_i = x_k$ or $l_i = \bar{x}_k$ s.t l_i is false on $v_i(k)$. Since $|C| = d+1$ and $|\{l_1, \dots, l_d\}| = d$ then by the Pigeonhole principle

there is at least one variable that occurs twice. w.l.o.g assume that l_1 and l_2 contain the same variable.

- If $l_1 = l_2$: then l_1 is false on $v_1(k)$ and l_2 is false on $v_2(k)$ and then $h_1(v_1) = h_2(v_2) = 0$ contradiction.
- If $l_1 = \neg l_2$: since $d \geq 2 \Rightarrow d+1 \geq 3$ then v_3 exists and either l_1 or l_2 is false on $v_3 \Rightarrow$ either $h_1(v_3) = 0$ or $h_2(v_3) = 0$ contradiction.

Hence, $d \leq VCdim(\mathcal{H}) \leq d \Rightarrow VCdim(\mathcal{H}) = d$.

4.Solution:

Let $h_S = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \{L_S(h)\}$ and $h^* = \underset{h \in \mathcal{H}}{\operatorname{min}} L_{\mathcal{D}}(h)$.

Claim : For every distribution \mathcal{D} and sample $S \sim D^m$ it holds that:

$$L_{\mathcal{D}}(h_S) - h^* \leq 2 \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$$

Proof:

$$\begin{aligned} L_{\mathcal{D}}(h_S) - h^* &= L_{\mathcal{D}}(h_S) - L_S(h_S) + L_S(h_S) - h^* \\ &\leq L_{\mathcal{D}}(h_S) - L_S(h_S) + L_S(h^*) - h^* \leq 2 \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \end{aligned}$$

Back to the question : Let $\epsilon, \delta > 0$. By the claim we get:

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2} \right] = \mathbb{P}_{S \sim \mathcal{D}^m} \left[2 \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \epsilon \right] \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D}}(h_S) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \leq \epsilon \right]$$

And if $\forall h \in \mathcal{H} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2}$ then $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2}$, so we get:

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\forall h \in \mathcal{H} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2} \right] \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2} \right] \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D}}(h_S) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \leq \epsilon \right]$$

Then by the uniform convergence property $\forall m \geq m(\frac{\epsilon}{2}, \delta)$ it holds that:

$$1 - \delta \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left[\forall h \in \mathcal{H} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{\epsilon}{2} \right] \leq \mathbb{P}_{S \sim \mathcal{D}^m} \left[L_{\mathcal{D}}(h_S) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) \leq \epsilon \right]$$

7.Solution:

Let $\mathcal{H}^1 \subseteq \mathcal{H}^2$ be two classes for binary classification and let $C \subset \mathcal{X}$ s.t $C = \{c_1, \dots, c_m\}$.

By definition $\mathcal{H}_C^1 = \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}^1\}$ and since $\mathcal{H}^1 \subseteq \mathcal{H}^2$ it holds :

$$\begin{aligned} \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}^1\} &\subseteq \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}^1\} \cup \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}^2 \setminus \mathcal{H}^1\} \\ &= \{h(c_1), \dots, h(c_m) : h \in \mathcal{H}^2\} \end{aligned}$$

Thus, $\mathcal{H}_C^1 \subseteq \mathcal{H}_C^2$. Therefore if \mathcal{H}^1 shatters C then

$$|\mathcal{H}_C^1| = 2^{|C|} \leq |\mathcal{H}_C^2| \leq 2^{|C|} \Rightarrow |\mathcal{H}_C^2| = 2^{|C|}$$

$\Rightarrow \mathcal{H}_C^2$ shatters C . Now, by the definition of $VCdim$ we get :

$$VCdim(\mathcal{H}^1) \leq VCdim(\mathcal{H}^2)$$

8.1.Solution:

$\tau(m)$ is the maximum number of ways that m points can be classified by using \mathcal{H} .

8.2.Solutions:

If $VCdim(\mathcal{H}) = \infty$ then for all $m \in \mathbb{N}$ there exists C with $|C| = m$ s.t $|\mathcal{H}_C| = 2^m \Rightarrow$ for all m it holds that $\tau(m) = 2^m$

8.3.Solution:

If $m \leq d$ then there exists C with $|C| = m$ s.t $|\mathcal{H}_C| = 2^m$ and since $\forall C \subset \mathcal{X}$ it holds that $|\mathcal{H}_C| \leq 2^{|C|}$ then $\tau(m) = 2^m$.

8.4.1.Solution:

Claim : for any finite $C \subset \mathcal{X}$

$$|\mathcal{H}_C| \leq |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}|$$

Proof : Let $C = \{c_1, \dots, c_m\}$. We proof the claim by induction on m .

Base : ($m = 1$) Let $B \subseteq C$ then $B = C = \{c_1\}$ or $B = \emptyset \Rightarrow |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}| \leq 2$.

- We know that $\mathcal{H} \text{ shutters } \emptyset \Rightarrow 1 \leq |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}| \leq 2$
- If $|\mathcal{H}_C| \leq 1$ then $|\mathcal{H}_C| \leq |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}|$
- Otherwise $|\mathcal{H}_C| = |\{h(c_1) : h \in \mathcal{H}\}| = 2 = 2^{|C|} \Rightarrow \mathcal{H} \text{ shutters } \emptyset \text{ and } \{c_1\} \Rightarrow |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}| = 2$
and we are done

Step : Assume the statment holds for any $k < m$ and $1 < m$:

Denote $C' = \{c_2, \dots, c_m\}$ and define $\mathcal{H}^C = \{h \in \mathcal{H} : \exists h' \in \mathcal{H} \text{ s.t } h|_{C'} \equiv h'|_{C'} \text{ and } h(c_1) \neq h'(c_1)\}$. We claim that $|\mathcal{H}_C| = |\mathcal{H}_{C'}| + |\mathcal{H}_{C'}^C|$:

For any $\bar{y} = (y_2, \dots, y_m)$ and $i \in \{0, 1\}$ define :

$$I_i(\bar{y}) = \begin{cases} 1 & (i, y_2, \dots, y_m) \in \mathcal{H}_C \\ 0 & \text{else} \end{cases}$$

And, denote $X_k = \{\bar{y} \in \mathcal{H}_{C'} : I_0(\bar{y}) + I_1(\bar{y}) = k\}$ for $k \in \{1, 2\}$. Now clearly $|\mathcal{H}_C| = |X_1| + |X_2|$.

Moreover, since $(y_1, y_2, \dots, y_m) \in \mathcal{H}_C \Leftrightarrow (y_2, \dots, y_m) \in \mathcal{H}_{C'}$ then $|X_0| = |\mathcal{H}_{C'}|$,

and $\bar{y} \in X_2 \Leftrightarrow (1, y_2, \dots, y_m), (0, y_2, \dots, y_m) \in \mathcal{H}_C \Leftrightarrow \bar{y} \in \mathcal{H}_{C'}^C$. Therefore, $|\mathcal{H}_C| = |\mathcal{H}_{C'}| + |\mathcal{H}_{C'}^C|$.

Let us examine the two sets:

- $\mathcal{H}_{C'}$: By the I.H it holds that $|\mathcal{H}_{C'}| \leq |\{B \subseteq C' : \mathcal{H} \text{ shutters } B\}|$
- $\mathcal{H}_{C'}^C$: By the I.H it holds that $|\mathcal{H}_{C'}^C| \leq |\{B \subseteq C' : \mathcal{H}^C \text{ shutters } B\}|$.

Let $B = \{b_1, \dots, b_k\} \subseteq C'$.

If \mathcal{H}^C shutters B then by the definition of $\mathcal{H}^C \forall h \in \mathcal{H}^C$ if $(h(b_1), \dots, h(b_k)) \in \mathcal{H}_B^C$ then :

$$(1, h(b_1), \dots, h(b_k)), (0, h(b_1), \dots, h(b_k)) \in \mathcal{H}_{B \cup \{c_1\}}^C$$

Hence $|\mathcal{H}_B^C| = 2^k \Rightarrow |\mathcal{H}_{B \cup \{c_1\}}^C| = 2^{k+1} \Rightarrow \{B \subseteq C' : \mathcal{H}^C \text{ shutters } B\} \subseteq \{B \subseteq C' : \mathcal{H}^C \text{ shutters } B \cup \{c_1\}\},$

and thus :

$$|\mathcal{H}_{C'}^C| \leq \left| \{B \subseteq C' : \mathcal{H}^C \text{ shutters } B\} \right| \leq \left| \{B \subseteq C' : \mathcal{H}^C \text{ shutters } B \cup \{c_1\}\} \right|$$

Moreover, since $\mathcal{H}^C \subseteq \mathcal{H}$ then $\{B \subseteq C' : \mathcal{H}^C \text{ shutters } B \cup \{c_1\}\} \subseteq \{B \subseteq C' : \mathcal{H} \text{ shutters } B \cup \{c_1\}\}.$

Finally, $\{B \subseteq C' : \mathcal{H} \text{ shutters } B\} = \{B \subseteq C : \mathcal{H} \text{ shutters } B \text{ and } c_1 \notin B\}$

and $\{B \subseteq C' : \mathcal{H}^C \text{ shutters } B \cup \{c_1\}\} = \{B \subseteq C : \mathcal{H}^C \text{ shutters } B \text{ and } c_1 \in B\}$ we get :

$$|\mathcal{H}_C| = |\mathcal{H}_{C'}| + |\mathcal{H}_{C'}^C|$$

$$\leq |\{B \subseteq C : \mathcal{H} \text{ shutters } B \text{ and } c_1 \notin B\}| + |\{B \subseteq C : \mathcal{H} \text{ shutters } B \text{ and } c_1 \in B\}| = |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}|$$

8.4.2.Solution:

The number of possible labeling of points in C is not more then the number of subsets of C that there points have all possible labeling.

8.4.3.Solution:

Assume $d < m$ and $VCdim(\mathcal{H}) = d.$

Let $B \in \{B \subseteq C : \mathcal{H} \text{ shutters } B\}$ then $VC(\mathcal{H}) = d \Rightarrow |B| \leq d \Rightarrow |\{B \subseteq C : \mathcal{H} \text{ shutters } B\}|.$ Since for fixed $0 \leq k \leq d$ the number subsets of size k of a set of size m is $\binom{m}{k}$ then :

$$|\{B \subseteq C : \mathcal{H} \text{ shutters } B\}| \leq \sum_{k=0}^d \binom{m}{k}$$

8.4.4.Solution:

Let $d, m \in \mathbb{N}$ s.t $VCdim(\mathcal{H}) = d$, $d < m$ and let $C^* \subset \mathcal{X}$ s.t $\tau(m) = |\mathcal{H}_{C^*}|$ then according to (8.4.1) it holds that :

$$\tau(m) = |\mathcal{H}_{C^*}| \leq |\{B \subseteq C^* : \mathcal{H} \text{ shutters } B\}| \leq \sum_{k=0}^d \binom{m}{k} \leq \left(\frac{em}{d}\right)^d$$

8.4.5.Solution:

For non empty \mathcal{H} it holds that $\forall C \subset \mathcal{X}$ it holds that $|\mathcal{H}_C| \leq 2^{|C|}$ then If $m = d$ then there exists $C^* \subset \mathcal{X}$ s.t \mathcal{H} shatters $C^* \Rightarrow \tau(m) = |\mathcal{H}_{C^*}| = 2^m$. Therefore :

$$\tau(m) = 2^m \leq \sum_{k=0}^m \binom{m}{k} = 2^m < e^m$$

Hence, the inequality holds but it is not tight.

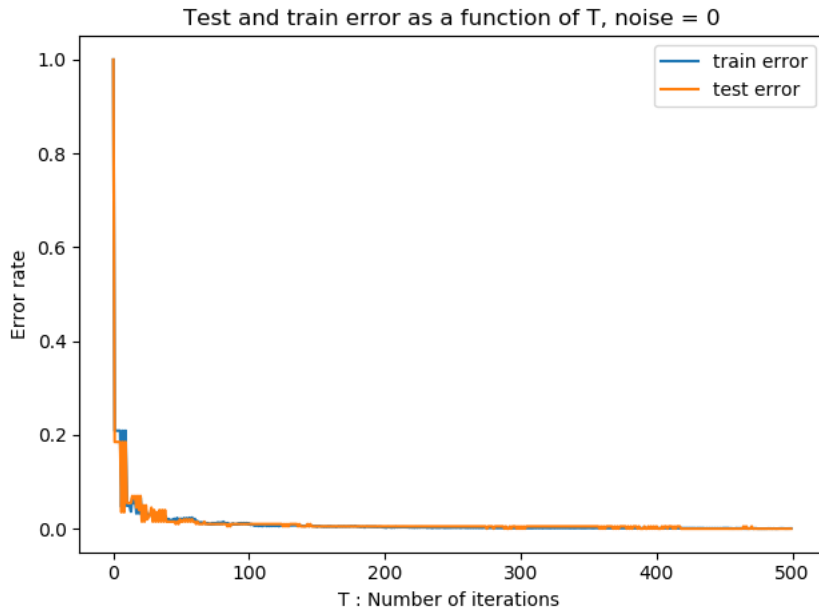
8.4.6.Solution:

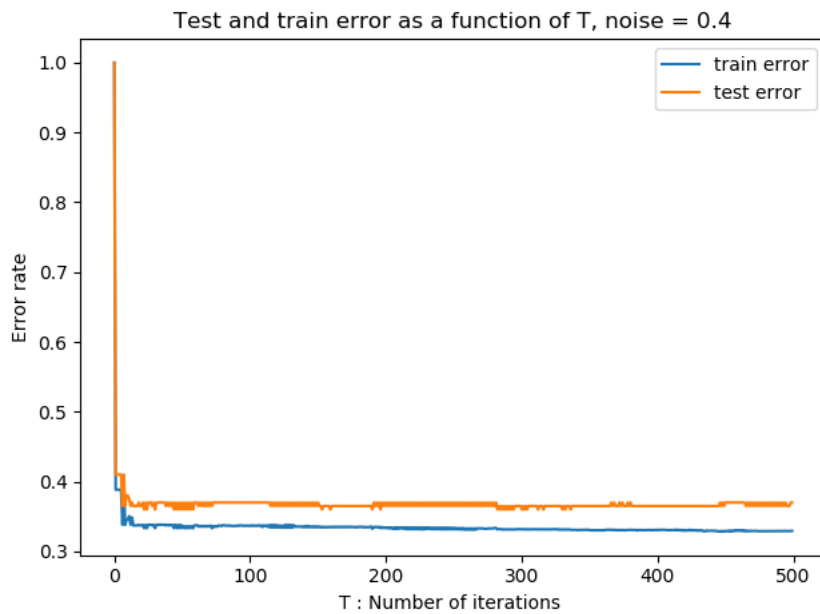
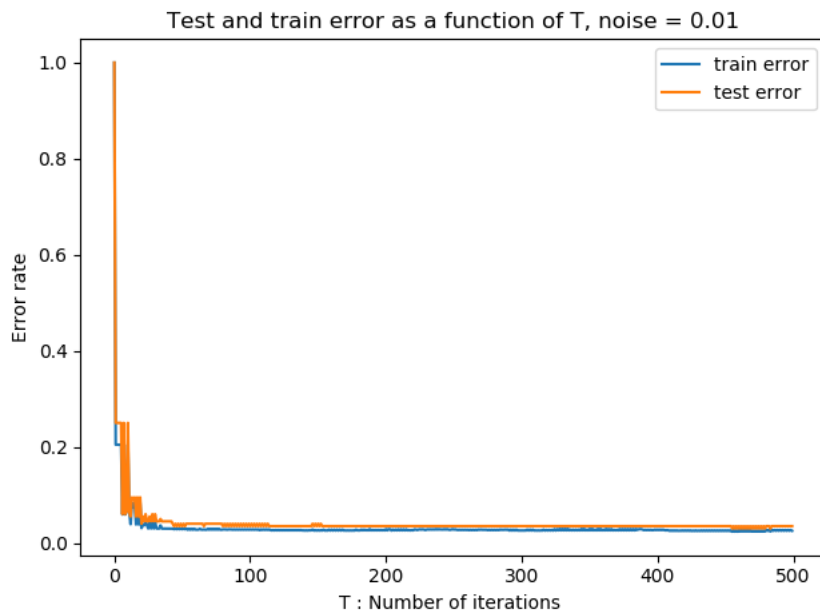
For $VCdim(\mathcal{H}) = d$.

- If $d < m$ then $\tau(m)$ grows polynomially in m : $\tau(m) \leq \left(\frac{em}{d}\right)^d$
- If $m \leq d$ then $\tau(m)$ grows exponentially in m : $\tau(m) = 2^m$

Now we can characterize $VCdim(\mathcal{H}) = d = \max\{m : \tau(m) = 2^m\}$ which means there is shattered sample of size d and there is no shattered sample of size greater than d .

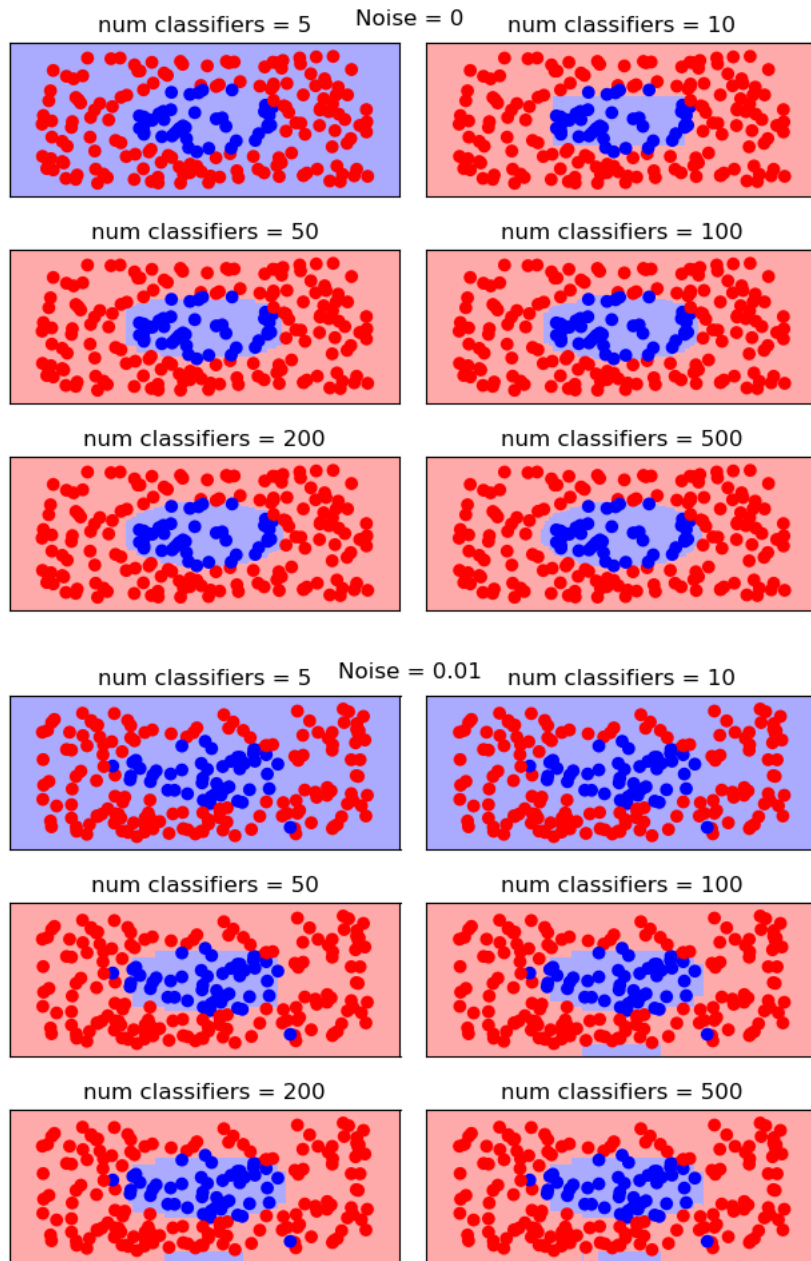
10.Solution:

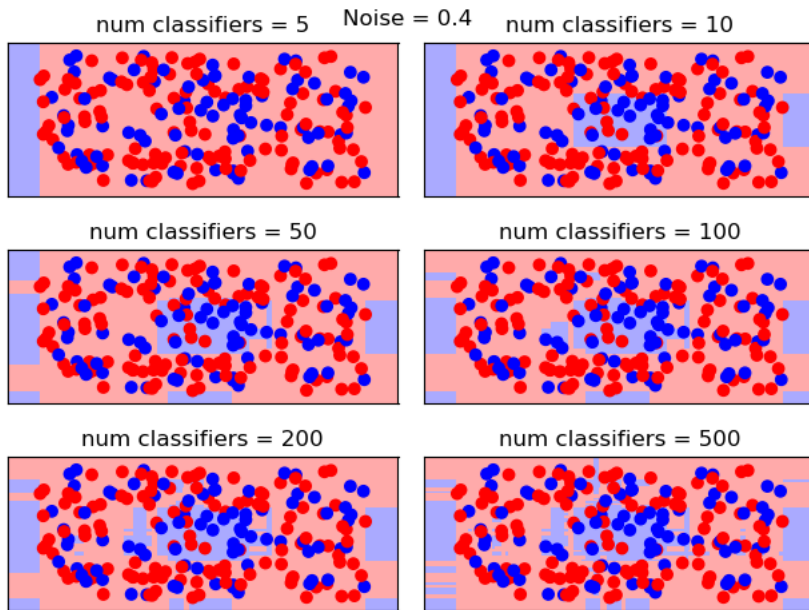




By looking at the graphs we can see that for data without noise both the bias and the variance are small and after a drastic change for small values of T there is only slight change as T grows (namely for T greater or equal to 200). Moreover, we can see that AdaBoost is sensitive to label noise. when the noise grows to 0.01 we can see a slight increase in the variance, but when the noise grows to 0.4 we can see clearly an increase in the bias and the variance grows with T . It seems that AdaBoost is sensitive to label noise. An explanation for this might be that AdaBoost is fitting a classification model to an exponential loss function, so every one mislabeled point could have a very strong influence on the final model learned.

11.Solution:



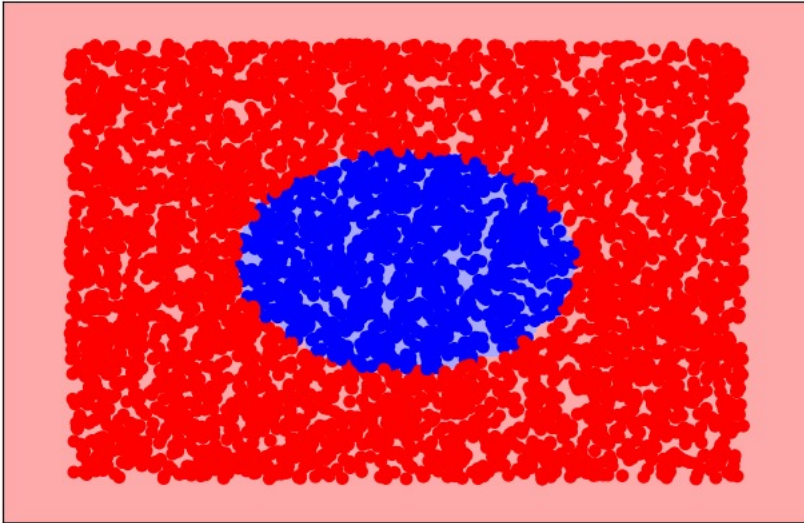


Here we can see a demonstration to the previous explanation. For not noisy data we can see that there are great results for only 100 classifiers. For noise of 0.01 we can see an outlier and it's effect on the final model learned.

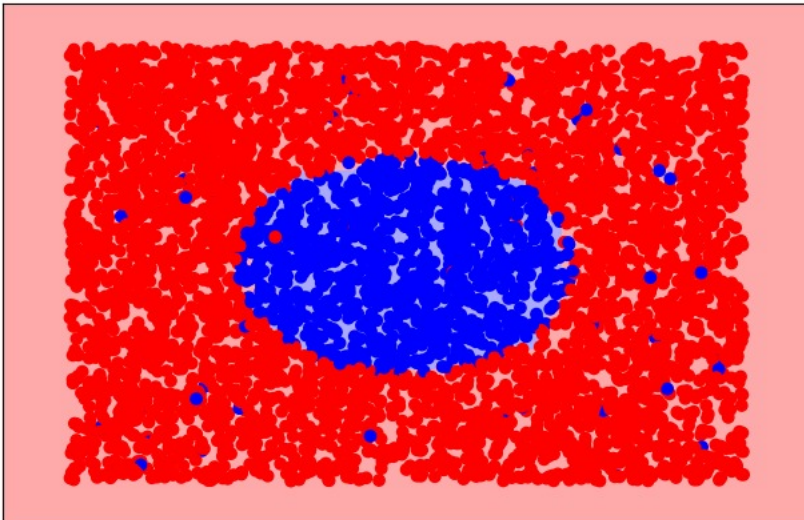
there is great improvment when the number of classifiers is 50 and then it's getting worse and seems to overfit and for noise of 0.4 we see very poor performance.

12.Solution:

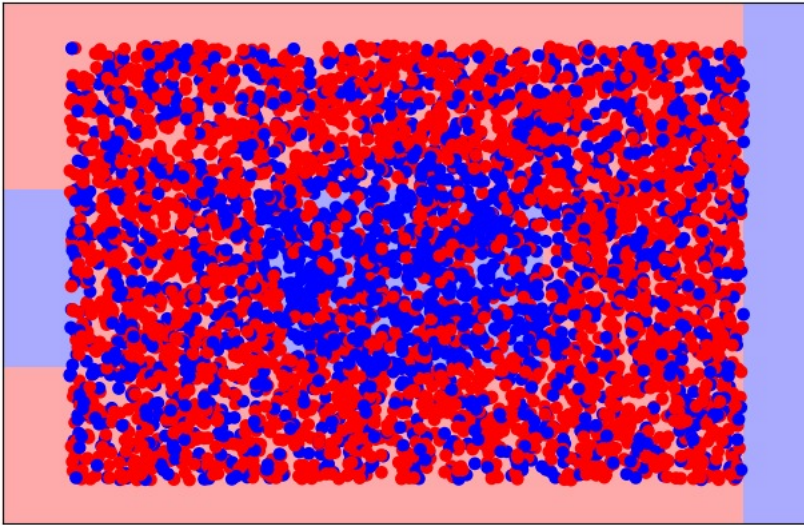
Noise = 0
Test Error: 0.005
T with minimal error: 500



Noise = 0.01
Test Error: 0.035
T with minimal error: 200



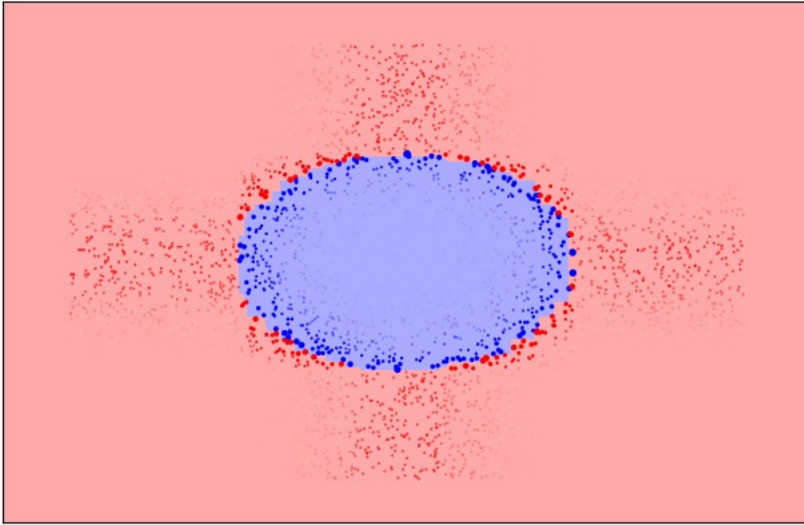
Noise = 0.4
Test Error: 0.31
T with minimal error: 10



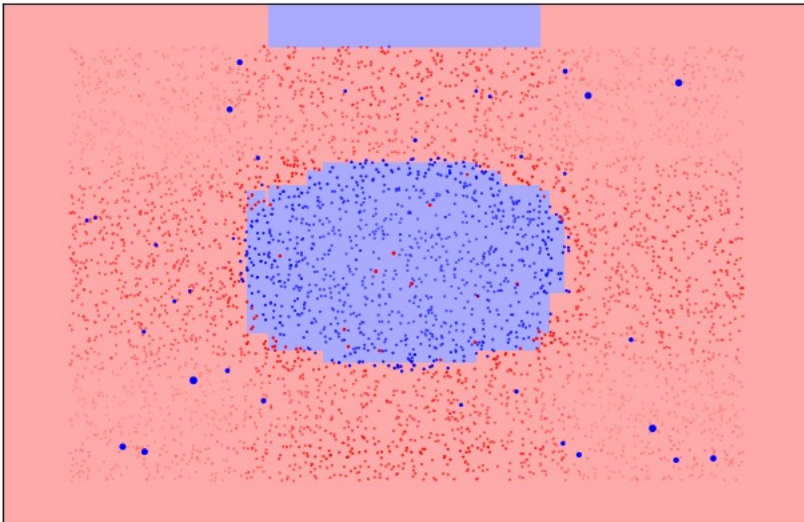
After running the algorithm few times it seems that typically when the noise grows the number of classifiers needed in order to gain minimum error decreasing. Again since every one mislabeled point could have a very strong influence. Since it is an additive model then at some point,for noisy date, it seems better to redunce the number of classifiers in order to reduce influence of the outliers.

13.Solution:

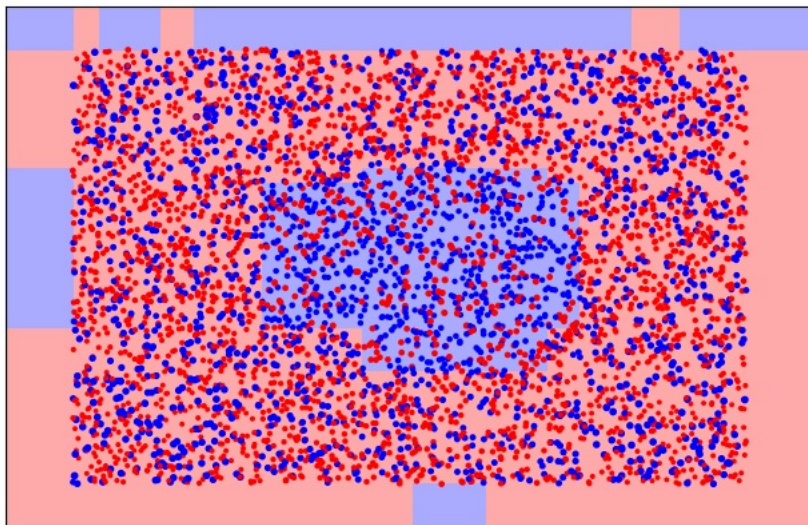
Noise = 0
Test Error: 0.005
T with minimal error: 500



Noise = 0.01
Test Error: 0.025
T with minimal error: 100



Noise = 0.4
Test Error: 0.315
T with minimal error: 50



We can see that for not noisy data the points around the circle get greater values which makes sense, but since there are no outliers the increasing of the weights around the circle seems to help and to get very well results. On the other hand the interesting part is a slight increase in noise (0.01), where only few outliers with large negative margin have a strong effect and it can be seen that the surrounding suffers more from misclassification. Finally for noise of 0.4 the data is too noisy and too many points have large values of weights.