

Seminar 2

① The Bayes optimal predictor has the smallest error among all possible classifiers

Let D be a probability distribution over $X \times \{0,1\}$.

The Bayes classifier is defined as:

$$f_D : X \rightarrow \{0,1\}, f_D(x) = \begin{cases} 1, & \underbrace{P_{(x,y) \sim D}(y=1|x)}_{\eta(x)} \geq \frac{1}{2} \\ 0, & \text{otherwise} \end{cases}$$

Let $g : X \rightarrow \{0,1\}$ be a random classifier. We want to show that $L_D(f_D) \leq L_D(g)$.

$$L_D(f_D) = E_{(x,y) \sim D} (\ell(f_D(x), y)) = E_{(x,y) \sim D} (1_{\{f_D(x) \neq y\}}) = P_{(x,y) \sim D} (f_D(x) \neq y)$$

$$\ell(f_D(x), y) = \begin{cases} 1, & f_D(x) \neq y \\ 0, & f_D(x) = y \end{cases} = 1_{\{f_D(x) \neq y\}}$$

$$E_{(x,y) \sim D} (1_{\{f_D(x) \neq y\}}) = E_{x \sim D} \left[E_{y \sim D|x} \left[1_{\{f_D(x) \neq y\}} | x \right] \right]$$

$$= P_{y \sim D|x} (f_D(x) \neq y)$$

$$P_{y \sim D|x} (f_D(x) \neq y | x) = P(y=1 | x) * 1_{[\eta(x) < \frac{1}{2}]} + P(y=0 | x) * 1_{[\eta(x) \geq \frac{1}{2}]}$$

$$= \eta(x) * 1_{[\eta(x) < \frac{1}{2}]} + (1 - \eta(x)) * 1_{[\eta(x) \geq \frac{1}{2}]}$$

$$= \begin{cases} \eta(x), & \eta(x) < \frac{1}{2} \\ 1 - \eta(x), & \eta(x) \geq \frac{1}{2} \end{cases} = \min(\eta(x), 1 - \eta(x))$$

$$P_D(g) = E_{(x,y) \sim D} (\ell(g(x), y)) = E_{(x,y) \sim D} (1_{g(x) \neq y}) = E_{x \sim D} \left[E_{y \sim D|x} \left[1_{\{g(x) \neq y\}} | x \right] \right]$$

$$P_{y \sim D|x} (g(x) \neq y | x) = P(g(x)=0, y=1 | x) + P(g(x)=1, y=0 | x)$$

$$\approx P(g(x)=0 | x) * P(y=1 | x) + P(g(x)=1 | x) * P(y=0 | x)$$

$$= P(g(x)=0 | x) * \underbrace{\eta(x)}_{\geq \min(\eta(x), 1 - \eta(x))} + P(g(x)=1 | x) * \underbrace{(1 - \eta(x))}_{\geq \min(\eta(x), 1 - \eta(x))}$$

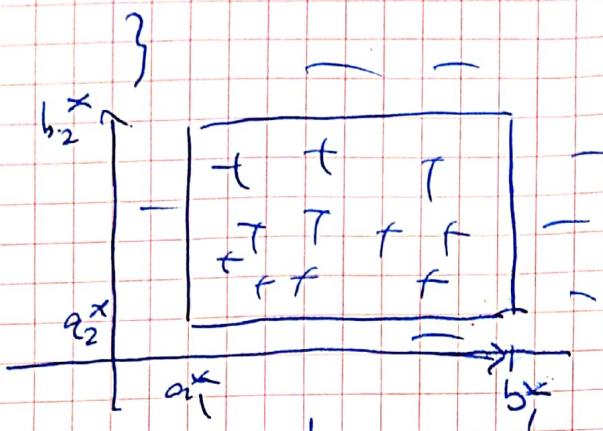
$$\geq \min(\eta(x), 1 - \eta(x)) + \min(\eta(x), 1 - \eta(x))$$

$$\geq (P(g(x)=0 | x) + P(g(x)=1 | x)) * \min(\eta(x), 1 - \eta(x))$$

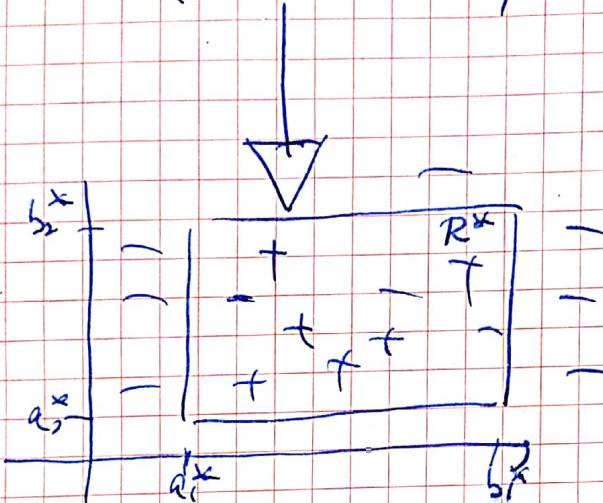
$$= \min(\eta(x), 1 - \eta(x)) = P_{y \sim D|x} (f_D(x) \neq y | x)$$

$$2) \quad \mathcal{H}_{\text{rect}} = \{h_{(a_1, b_1, a_2, b_2)} : \mathbb{R}^2 \rightarrow \{0, 1\}, a_1 \leq b_1, a_2 \leq b_2\}$$

$$h_{(a_1, b_1, a_2, b_2)}(x_1, x_2) = \begin{cases} 1, & x \in [a_1, b_1] \times [a_2, b_2] \\ 0, & \text{otherwise} \end{cases}$$



PAC-learnability
in scenario 1



agnostic PAC
learnability

Positive labels are
flipped with probability
 $0 < \eta < \frac{1}{2}$.

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

$$y_i \equiv \begin{cases} 0, & \text{if } x_i \notin [a_1^*, b_1^*] \times [a_2^*, b_2^*] \\ 0, & \text{with prob. } \eta \text{ if } x_i \in [a_1^*, b_1^*] \times [a_2^*, b_2^*] \\ 1, & \text{with prob. } 1-\eta \text{ if } x_i \in [a_1^*, b_1^*] \times [a_2^*, b_2^*] \end{cases}$$

$$R^* = [a_1^*, b_1^*] \times [a_2^*, b_2^*]$$

The chance to get a training point labeled as positive is to sample a point from R^* and the label is not flipped so the chance is

$$\Pr(R^* \cap (1-\eta))$$

A - the learning algorithm that returns the tightest rectangle containing positive points.

$$h_S = A(S), \quad h_S = h_{(a_{1S}, b_{1S}, a_{2S}, b_{2S})}$$

$$a_{1S} = \min_{(x_{i,1}) \in S} x_{i,1}$$

$$a_{2S} = \min_{(x_{i,1}) \in S} x_{i,2}$$

$$b_{1S} = \max_{(x_{i,1}) \in S} x_{i,1}$$

$$b_{2S} = \max_{(x_{i,1}) \in S} x_{i,2}$$

If S doesn't contain positive samples then return let $z = (z_1, z_2)$
 and $h_S = h_{(z_1, z_2, z_1, z_2)}$, where z is a random point in \mathbb{R}^2

We want to show that H_{rec}^2 is agnostic PAC-learnable;

there exist a function $m_{H_{\text{rec}}^2}: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A such that:

for every $\epsilon > 0$, for every $\delta > 0$, for every distribution D over $Z = \mathbb{R}^2 \times \{0, 1\}$, $D = D_X \times D_Y$

when we run the learning algorithm A on a training set S consisting of $n \geq m_{H_{\text{rec}}^2}(\epsilon, \delta)$ examples sampled i.i.d from D the algorithm A return a hypothesis $h_S = A(S)$ from H_{rec}^2 such that

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

In our case the smallest achievable real error is

$$\min_{h \in H} L_D(h) = L_D(h^*) = \eta \cdot D_X(R^*)$$

Consider $\epsilon > 0$, $\delta > 0$ and D_X a distribution over \mathbb{R}^2 .

Case 1: if $D_X(R^*) \leq \epsilon \Rightarrow h_S$ can only make errors on points inside R^*

$$\therefore \Pr_{S \sim D^m} (L_D(h_S) \leq \epsilon) = 1 \quad v.$$

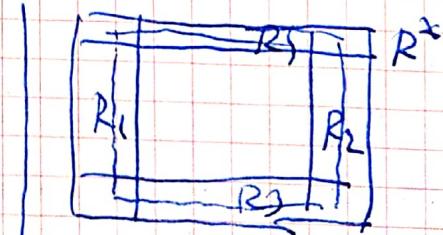
Case 2: if $D_X(R^*) > \epsilon$

Construct rectangles R_1, R_2, R_3, R_4 (here in similar)

$$\text{such that } D_X(R_i) = \frac{\epsilon}{4}$$

i) if $h_S = A(S)$ intersects all R_i , ($i \geq 1$)
 then h_S will make errors on:

$$P_S > \text{because of flipping} \rightarrow D(R_S) = \eta \rightarrow R^* \setminus R \rightarrow D_X(R^* \setminus R) \leq \epsilon$$



So in this case we have:

$$\Pr_{\text{random}}(L_D(h_S) \leq \gamma \cdot D(R^{\star} + \varepsilon)) = 1.$$

ii) If $h_S = A(S)$ doesn't intersect a rectangle R_i

Disk $P_i = \{S | S \cap D_x^m \neq \emptyset \text{ s.t. } R_S - \text{the rectangle formed by } A(S) \text{ doesn't intersect } P_i\}$

$$\Pr_{S \sim D^m} (L_S(h_S) > \eta \cdot \mathcal{D}(P^*) + \varepsilon) \leq \sum_i \mathcal{D}_x(F_i) \text{ sehr.}$$

$Q^m(F_i)$ = the probability of sampling m points and none of them

If a positive point in R

$$= \left(1 - \frac{\xi}{\zeta} + \frac{\xi \cdot \eta}{\zeta} \right)^m = \left(1 - \frac{\xi}{\zeta} (1 - \eta) \right)^m$$

bracket below $\frac{\xi}{\zeta}$: probability of sampling
 bracket below $\frac{\xi \cdot \eta}{\zeta}$: part of sample
 under ξ : part in P_{ij}
 under η : part outside P_{ij} but staying its label

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

$$1 - \frac{\xi}{\eta} (1-\eta) \leq e^{-\frac{\xi}{\eta} (1-\eta)}$$

$$\text{So: } P_{S=2^m} (L_S(h_S) \geq \gamma \cdot D(R^*) / \epsilon) \leq \frac{\epsilon^{-\frac{\epsilon}{4}(1-\gamma)}}{2^m} < \frac{\delta}{2}$$

$$G = \{e^{-\frac{2\pi i}{n}(r-n)/m}\}$$

$$e^{-\frac{\epsilon}{3}(1-n)^m} \leq \frac{d}{n} \quad | \log$$

$$m \cdot \left(-\frac{\varepsilon}{\zeta}\right)(1-\eta) < \log \frac{\sigma}{\eta} \quad \left| - \left(-\frac{\varepsilon}{\zeta}\right) \cdot \frac{1}{1-\eta}\right.$$

$$m > -\frac{4}{\varepsilon} \cdot \frac{1}{1-\eta} \cdot \log \frac{\sigma}{\delta}$$

$$m > \frac{4}{\varepsilon} \frac{1}{1-\eta} \cdot \log \frac{1}{\delta}$$

$$(3) \quad \mathcal{C} = \mathcal{H}_{\text{intervals}} = \{ h_{a,b} : \mathbb{R} \rightarrow \{0,1\}, a \leq b \}$$

$$h_{a,b}(x) = \mathbb{1}_{[a,b]}(x) = \begin{cases} 1, & x \in [a,b] \\ 0, & \text{otherwise} \end{cases}$$

Consider a training set $S_2 \{ (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m) \}$

We are in the realizable case, $\exists h^* \in \mathcal{H}_{\text{intervals}}$ that labels the examples, $y_i = h^*(x_i)$



We want to show that $\mathcal{H}_{\text{intervals}}$ is PAC-learnable.

Consider A the learning algorithm that gets sample S and outputs $h_S = A(S) = \text{the tightest interval containing all the positive examples.}$

$$h_S \sim h_{a_S, b_S}, \quad a_S = \min_{(x_i, 1) \in S} x_i, \quad b_S = \max_{(x_i, 1) \in S} x_i, \quad R_S = [a_S, b_S]$$

If there is no $(x_i, 1) \in S$ (S doesn't contain positive examples) then $a_S = b_S = 2$ a random point ~~not~~ such that $(2, 0) \notin S$.

From construction we see that $L_{D, h^*}(h_S) = 0$.

Let $\varepsilon > 0$, $\delta > 0$ and D a distribution over \mathbb{R} : What should be $m \geq m(\varepsilon, \delta)$

$$\text{such that } \Pr_{S \sim D^m} (L_{D, h^*}(h_S) > \varepsilon) < \delta$$

Case 1: If $D([a^*, b^*]) \leq \varepsilon$ then $\Pr_{S \sim D^m} (L_{D, h^*}(h_S) > \varepsilon) = 0 \checkmark$

Case 2: If $D([a^*, b^*]) > \varepsilon$

Consider R_1 and R_2 , $R_1 = [a_1^*, a_1]$, $R_2 = [b_2, b_2^*]$

such that $D(R_1) = D(R_2) = \frac{\varepsilon}{2}$.

If $R_S \cap R_1 \neq \emptyset$ and $R_S \cap R_2 \neq \emptyset$ then $\Pr_{S \sim D^m} (L_{D, h^*}(h_S) > \varepsilon) = 0 \checkmark$

$$\text{Otherwise } \Pr_{S \sim D^m} (L_{D, h^*}(h_S) > \varepsilon) \leq 2 \cdot \left(1 - \frac{\varepsilon}{2}\right)^m \leq 2 \cdot e^{-\frac{\varepsilon^2 m}{2}} < \delta$$

$$\Rightarrow \boxed{m > \frac{2}{\varepsilon^2} \log \frac{2}{\delta}}$$

④ PAC-learning algorithm for the class \mathcal{C}_2 formed by unions of two closed intervals : $\mathcal{C}_2 = \{ h_{(a,b,c,d)} : \mathbb{R} \rightarrow [0,1], h_{(a,b,c,d)} = \mathbb{1}_{[a,b] \cup [c,d]} \}$

$$a \leq b \leq c \leq d, h_{(a,b,c,d)}(x) = \begin{cases} 1, & x \in [a,b] \cup [c,d] \\ 0, & \text{otherwise} \end{cases}$$

Consider $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$, ~~where $y_i \in \{0, 1\}$~~

where $y_i = h^*(x_i)$, $h^* = h_{(a^*, b^*, c^*, d^*)}$ - realizability assumption

$$\overline{\underset{a^*}{\left[\begin{array}{cccc} + & + & + & T \end{array} \right]}} = \overline{\underset{b^*}{\left[\begin{array}{cccc} + & + & + & \cdot \end{array} \right]}} = \overline{\underset{c^*}{\left[\begin{array}{cccc} + & + & \cdot & \cdot \end{array} \right]}} = \overline{\underset{d^*}{\left[\begin{array}{cccc} + & \cdot & \cdot & \cdot \end{array} \right]}}$$

Consider the following learning algorithm A that takes input S :

- Sort S in ascending order of x_i
- Go over the sorted training examples and take the intervals where consecutive training example labeled as positive start and end the interval.

You can obtain one or two intervals.

- If you obtained just one interval you can have $a_s = \min_{y_i=1} x_i$, $b_s = \max_{y_i=1} x_i$, $c_s = d_s = b_s$.

If you obtained two intervals then $a_s = \min_{y_i=1} x_i$ $a_s \leq b_s < c_s \leq d_s$

$$d_s = \max_{y_i=1} x_i$$

Return $h_S = h_{(a_s, b_s, c_s, d_s)} = \mathbb{1}_{[a_s, b_s] \cup [c_s, d_s]}$

We need to find $m \geq m_{\mathcal{C}_2}(\epsilon, \delta)$ such that for $\epsilon > 0$ and for every D distribution over \mathbb{R} we have that

$$\Pr_{S \sim D^m} (\mathcal{L}_{D, h^*}(h_S) > \epsilon) < \delta$$

Let $\epsilon > 0, \delta > 0$ and let D be a distribution over \mathbb{R} .

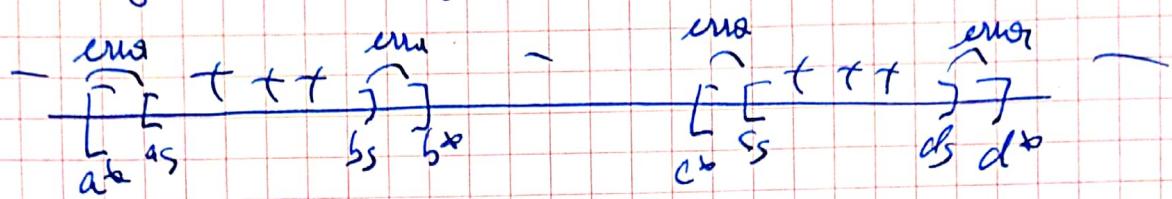
The region when h_S makes error is always $\subseteq [a^*, d^*]$

Case 1 So, if $D([a^*, d^*]) \leq \epsilon$ then $\Pr_{S \sim D^m} (\mathcal{L}_{D, h^*}(h_S) > \epsilon) = 0$.

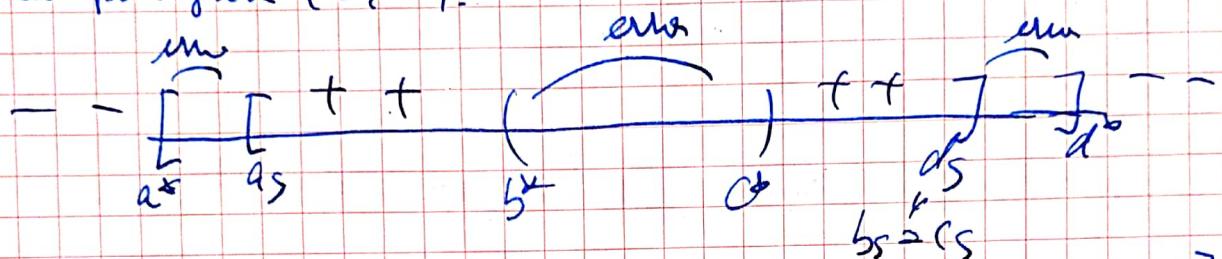
$$(\text{case 2}) \quad \mathcal{D}([a^*, b^*]) > \varepsilon$$

These types of error that h_S can make are:

- false negatives in $[a^*, b^*]$ and $[c^*, d^*]$



- false positive in (b^*, c^*) if sample S does not contain any point sampled from (b^*, c^*) .



$$\begin{aligned} \text{Define } L_{FP}(h_S) &= \underset{x \sim D}{P} [x \in [a_S^*, b_S^*] \setminus ([a^*, b^*] \cup [c^*, d^*])] \\ &= \underset{x \sim D}{P} [x \in (b^*, c^*) \subseteq [a_S^*, b_S^*] \cup [c_S^*, d_S^*]] \end{aligned}$$

$$L_{FN,1}(h_S) = \underset{x \sim D}{P} (x \in [a^*, b^*] \setminus [a_S^*, b_S^*])$$

$$L_{FN,2}(h_S) = \underset{x \sim D}{P} (x \in [c^*, d^*] \setminus [c_S^*, d_S^*])$$

So, if we want to have $L_{D,h_S}(\varepsilon) > \varepsilon$ then one of the numbers

$$L_{FP}(h_S), L_{FN,1}(h_S), L_{FN,2}(h_S) \text{ must be } \geq \frac{\varepsilon}{3}.$$

$$\text{Define } F_1 = \{S \sim D^m \mid L_{FP}(h_S) > \frac{\varepsilon}{3}\}$$

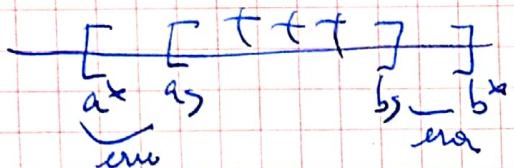
$$F_2 = \{S \sim D^m \mid L_{FN,1}(h_S) > \frac{\varepsilon}{3}\}$$

$$F_3 = \{S \sim D^m \mid L_{FN,2}(h_S) > \frac{\varepsilon}{3}\}$$

$$\text{So, } \underset{S \sim D^m}{P}(L_{D,h_S}(\varepsilon) > \varepsilon) \leq \underset{S \sim D^m}{P}(F_1 \cup F_2 \cup F_3) \leq \sum_{i=1}^3 P(F_i)$$

$P(F_1) = P_{\text{SDM}}(L_{FP}(h_S) > \frac{\epsilon}{3})$ (this means that $D((b^*, c^*)) > \frac{\epsilon}{3}$)
 and no point from (b^*, c^*) is sampled in S) $\leq (1 - \frac{\epsilon}{3})^m \leq e^{-\frac{\epsilon}{3}m}$

$$P(F_2) = P_{\text{SDM}}(L_{FN,L}(h_S) > \frac{\epsilon}{3})$$



Construct $R_1 = [a^*, a]$ and $R_2 = [b^*, b]$ such that

$$D(R_1) \geq D(R_2) \geq \frac{\epsilon}{6}.$$

if $[a_S, b_S] \cap R_1 \neq \emptyset$ and $[a_S, b_S] \cap R_2 \neq \emptyset$ then the user
 made by h_S on (a^*, b^*) is smaller than $\frac{\epsilon}{2} + \frac{\epsilon}{2} = \frac{\epsilon}{3}$.

$$\text{So } L_{FN,L}(h_S) > \frac{\epsilon}{3} \Rightarrow [a_S, b_S] \cap R_1 = \emptyset \text{ or } [a_S, b_S] \cap R_2 = \emptyset.$$

$$\text{define } F_{21} = \{ \text{SDM} | [a_S, b_S] \cap R_1 = \emptyset \}$$

$$F_{22} = \{ \text{SDM} | [a_S, b_S] \cap R_2 = \emptyset \}$$

$$P(F_2) \leq P(F_{21} \cup F_{22}) \leq P(F_{21}) + P(F_{22}) = 2 \cdot (1 - \frac{\epsilon}{6})^m \leq 2 \cdot e^{-\frac{\epsilon}{3}m}$$

In the same way we can prove that

$$P(F_3) \leq 2 \cdot e^{-\frac{\epsilon}{6}m}.$$

So we obtain that

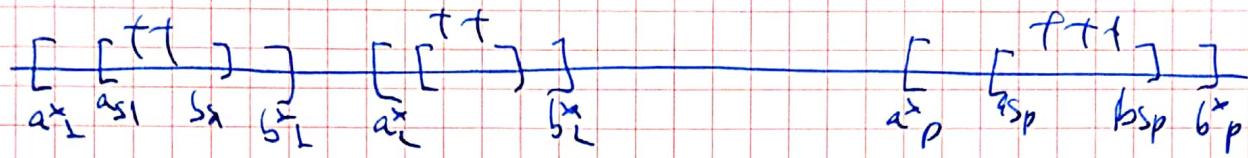
$$P_{\text{SDM}}(L_{\theta, h^*}(h_S) > \epsilon) \leq e^{-\frac{\epsilon}{3}m} + 4 \cdot e^{-\frac{\epsilon}{6}m} \leq e^{-\frac{\epsilon}{6}m} + 4e^{-\frac{\epsilon}{6}m} = 5e^{-\frac{\epsilon}{6}m} < \delta.$$

$$\Rightarrow e^{-\frac{\epsilon}{6}m} < \frac{\delta}{5} \quad |\lg \Rightarrow -\frac{\epsilon}{6}m < \lg \frac{\delta}{5} \quad | \cdot (-\frac{6}{\epsilon})$$

$$\boxed{m > \frac{6}{\epsilon} \cdot \lg \frac{5}{\delta}}$$

In the general case, for $\mathcal{C}_p = \text{union of } p \text{ intervals}$, the proof is similar, the only difference is that:

- there are $(p-1)$ regions of false positives
- $2p$ regions of false negatives



So we have

$$\left| m \geq \frac{2(2p-1)}{\epsilon} \times \log \frac{3+2p-1}{\delta} \right|$$

$$\left| m \geq \frac{2(2p-1)}{\epsilon} \times \log \frac{3p-1}{\delta} \right|$$

time complexity \rightarrow given by sorting $\leq O(m \log m)$

5) Let $h \in \mathcal{H}$, with $L_{(\bar{D}_m, g)}(h) > \varepsilon \Rightarrow$

$$\Pr_{x \sim \bar{D}_m} [h(x) \neq f(x)] > \varepsilon \Leftrightarrow \Pr_{x \sim \bar{D}_m} [h(x) = f(x)] = 1 - \Pr_{x \sim \bar{D}_m} [h(x) \neq f(x)] < 1 - \varepsilon$$

$\Pr_{x \sim \bar{D}_m} [h(x) = f(x)] = (\text{expect each point with } x \text{ can be sampled from each } D_i, \text{ with probability } \frac{1}{m}) =$

$$= \frac{1}{m} \cdot \Pr_{x \sim D_1} [h(x) = f(x)] + \dots + \frac{1}{m} \cdot \Pr_{x \sim D_m} [h(x) = f(x)] \\ = \frac{1}{m} \sum_{i=1}^m \Pr_{x \sim D_i} [h(x) = f(x)] < 1 - \varepsilon.$$

$$S = \{(x_1, f(x_1)), (x_2, f(x_2)), \dots, (x_m, f(x_m)) \mid \text{when } x_i \sim D_i\}$$

h consistent with S if $L_S(h) = 0$

$$\Pr_{S \sim D_1 \times D_2 \times \dots \times D_m} [L_S(h) = 0] = \prod_{i=1}^m \Pr_{x \sim D_i} [h(x_i) = f(x_i)]$$

$$= \prod_{i=1}^m \Pr_{x \sim D_i} [h(x) = f(x)]$$

$$= \left[\left(\prod_{i=1}^m \Pr_{x \sim D_i} [h(x) = f(x)] \right)^{\frac{1}{m}} \right]^m$$

= geometric mean = $(a_1 \cdot a_2 \cdot \dots \cdot a_m)^{\frac{1}{m}}$

When $a_i = \Pr_{x \sim D_i} [h(x) = f(x)]$ = probability that h correctly labels a point $x \sim D_i$

\leq arithmetic mean = $\left(\frac{a_1 + a_2 + \dots + a_m}{m} \right)$

$$\leq \left[\frac{1}{m} \sum_{i=1}^m \Pr_{x \sim D_i} [h(x) = f(x)] \right]^m < (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

There are at most $|\mathcal{H}|$ members of h hypothesis.

So, we obtain that

$$\Pr [\exists h \in \mathcal{H} \text{ s.t. } L_{(\bar{D}_m, g)}(h) > \varepsilon \text{ and } L_{(S, g)}(h) = 0] \leq |\mathcal{H}| \cdot e^{-\varepsilon m}$$