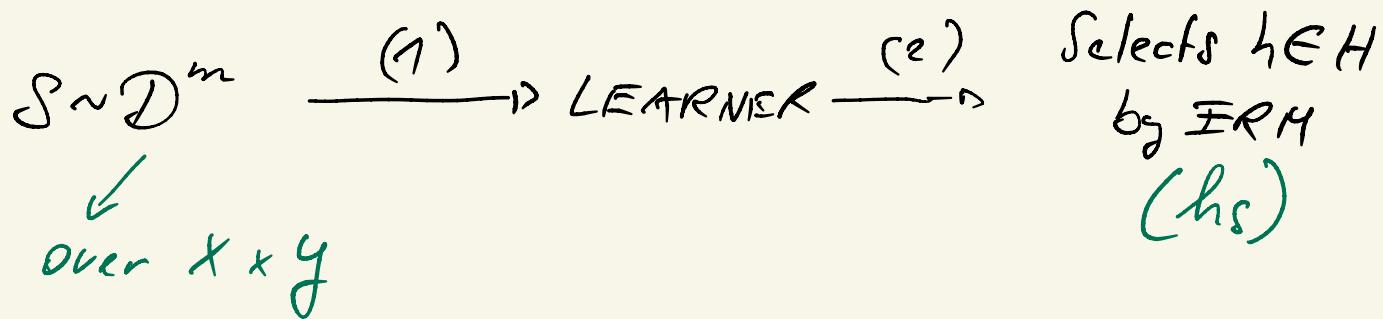


MACHINE LEARNING - VO 2020

KW III ROLAND



Uniform Convergence



Good property: Ideally $L_S(h_S)$ should be close to $L_D(h_S)$ — this should hold for any $h \in H$.

Def. A sample S is called ε -representative (with respect to D, H , loss function ℓ and $Z = X \times Y$) if

$$\forall h \in H: |L_S(h) - L_D(h)| \leq \varepsilon$$

Lemma: If a training sample, S , is $\frac{\varepsilon}{2}$ -representative, then any output of $\text{ERM}_H(S)$, say h_S , satisfies

$$L_D(h_S) \leq \min_{h \in H} L_D(h) + \varepsilon$$

$$\begin{aligned}
 \text{Proof: } L_D(h_S) &\leq L_S(h_S) + \frac{\varepsilon}{2} && \left(\text{by } \frac{\varepsilon}{2}\text{-rep.} \right) \\
 &\leq L_S(h) + \frac{\varepsilon}{2} && \left(\text{as } h_S \text{ is a } \hat{\text{ERH}} \text{ hyp.} \right) \\
 &\leq L_D(h) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} && \left(\text{by } \frac{\varepsilon}{2}\text{-rep.} \right) \\
 &= L_D(h) + \varepsilon
 \end{aligned}$$

Since this inequality chain holds for any $h \in H$, it also holds for the minimum.

$$\Rightarrow L_D(h_S) \leq \min_{h \in H} L_D(h) + \varepsilon \quad \boxed{P_2}$$

Def. A hypothesis class H has the uniform convergence property, if $\exists m_H^{uc}: (0, 1)^2 \rightarrow \mathbb{N}$, such that for every $\varepsilon, \delta \in (0, 1)$ and every prob. distribution D over $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$, if S is a sample of size $m \geq m_H^{uc}(0, \varepsilon)$, then, with prob. of at least $1 - \delta$, S is ε -representative.

We already know that finite h.yo. classes (H) are PAC learnable via ERM.

Question: Are finite h.yo. classes agnostic PAC learnable?

Let $S = (z_1, \dots, z_m)$, $z_i \sim D$, $z_i \in X \times Y$ (the z_i are drawn i.i.d from D).

$$D^m \left(\{S : \forall h \in H, |L_S(h) - L_D(h)| \leq \epsilon\} \right) \geq 1 - \delta$$

Equivalently,

$$D^m \left(\{S : \exists h \in H, |L_S(h) - L_D(h)| > \epsilon\} \right) \leq \delta$$

$$D^m \left(\bigcup_{h \in H} \{S : |L_S(h) - L_D(h)| > \epsilon\} \right) \stackrel{\text{Union bound}}{\leq}$$

$$(x) \sum_{h \in H} D^m \left(\{S : |L_S(h) - L_D(h)| > \epsilon\} \right) \quad \left| \begin{array}{l} L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)] \\ \epsilon = \frac{1}{m} \sum_{i=1}^m l(h, z_i) \end{array} \right.$$

Taking the expectation, gives us
 $L_D(h)$

If $l = l_{0-1}$ (0-1 loss), we know the lower and upper bound (0, 1).

By the Hoeffding inequality, we get

$$D^m \left(\{S : |L_S(h) - L_D(h)| > \epsilon\} \right) \leq 2 \cdot e^{-\epsilon^2 m / (b-a)^2} \quad \left| \begin{array}{l} b=1 \\ a=0 \\ (L_S - L_D) = 1 \end{array} \right.$$

$$= 2 \cdot e^{-\epsilon^2 m}$$

Overall (x), we get

$$D^m \left(\{S : \exists h \in H, |L_S(h) - L_D(h)| > \epsilon\} \right) \leq \sum_{h \in H} 2 \cdot e^{-\epsilon^2 m} = 2 \cdot |H| \cdot e^{-\epsilon^2 m}$$

For $2 \cdot |H| \cdot e^{-\epsilon^2 m}$ to be smaller than $\delta \in (0, 1)$, we need

$$m > \log\left(\frac{2 \cdot |H|}{\delta}\right) \cdot \frac{1}{2\epsilon^2}$$

Corollary: Let H be a finite hyp. class ($|H| < \infty$) and $\ell: H \times \mathcal{Z} \rightarrow [0, 1]$. Then H has the uniform convergence property with

$$m_H^{uc}(\epsilon, \delta) \leq \left\lceil \log\left(\frac{2 \cdot |H|}{\delta}\right) \cdot \frac{1}{2\epsilon^2} \right\rceil$$

Furthermore, H is agnostic PAC learnable by ERM with

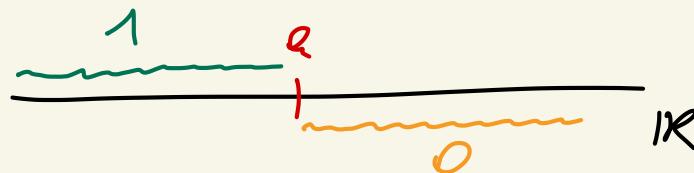
$$\underline{m_H(\epsilon, \delta)} \leq m_H^{uc}\left(\frac{\epsilon}{2}, \delta\right) \leq \left\lceil 2 \cdot \log\left(\frac{2 \cdot |H|}{\delta}\right) \cancel{\frac{1}{\epsilon^2}} \right\rceil$$

Last lecture: Finite H are agnostic PAC learnable.
(shown via uniform convergence).

Example (of an infinite hyp. class that is PAC learnable).

(*) $H = \{h_\alpha : \alpha \in \mathbb{R}\}$, $h_\alpha : \mathbb{R} \rightarrow \{0, 1\}$

$$h_\alpha(x) = \mathbb{1}_{[x < \alpha]} = \begin{cases} 1, & \text{if } x < \alpha \\ 0, & \text{else} \end{cases}$$



("class of threshold functions")

Lemma: Let H be the hyp. class (*) of threshold functions.
Then, H is PAC learnable ($\forall \epsilon \in \mathbb{R}_+$) with sample

complexity $m_H(\epsilon, \delta) \leq \left\lceil \frac{\log\left(\frac{2}{\delta}\right)}{\epsilon} \right\rceil$.

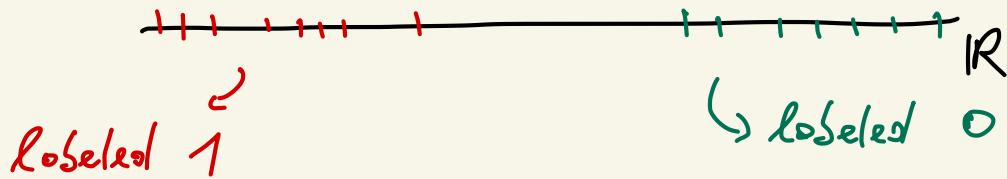
Proof: Let h^* be such that

$$L_{D, \epsilon}(h^*) = 0$$

and let $\alpha^* \in \mathbb{R}$ be the corresponding threshold α^* .

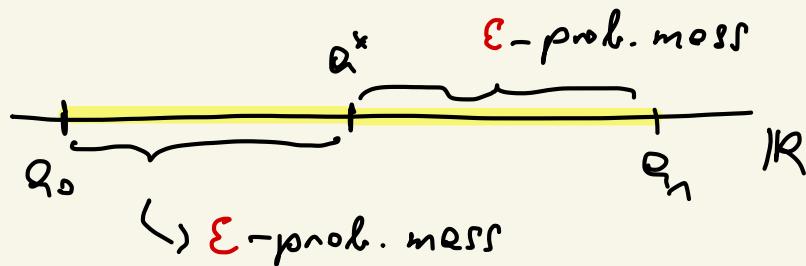


Let S be our training data.



$$b_0 = \max \{x : (x, 1) \in S\}$$

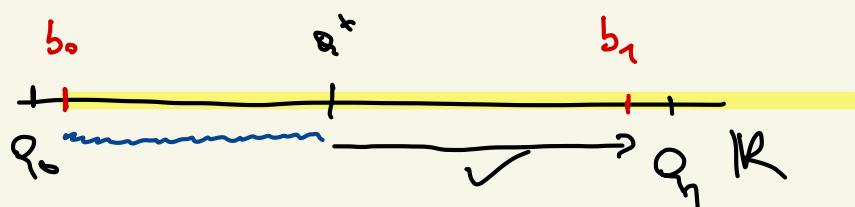
$$b_1 = \min \{x : (x, 0) \in S\}$$



We have $q_0 < q^* < q_n$ such that

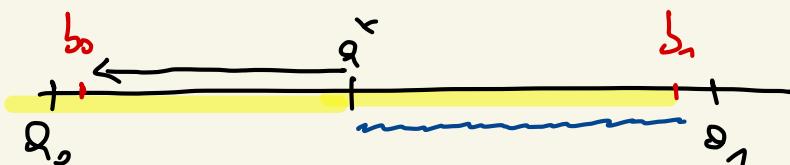
$$\mathcal{D}(\{x : x \in (q_0, q^*)\}) = \mathcal{D}(\{x : x \in (q^*, q_n)\}) = \epsilon$$

Case 1)



So b_0 is the thr. of the FPR hypothesis.

Case 2)



ERT picks b_1 as the thr.

\Rightarrow It is sufficient that $b_0 \geq Q_0$ AND $b_1 \leq Q_1$ for

$$L_{D,f}(h_S) \leq \varepsilon$$

with h_S being an ERH hypothesis. We know that the thresholds of h_S has to be within (b_0, b_1) , i.e., $b_S \in (b_0, b_1)$.

$$\begin{aligned} \Rightarrow \underset{S|x \sim D^m}{\mathbb{P}}[L_{D,f}(h_S) > \varepsilon] &\leq \underset{S|x \sim D^m}{\mathbb{P}} \left[b_0 < Q_0 \vee \underbrace{\left\{ b_1 > Q_1 \right\}}_{\text{OR}} \right] \\ &\leq \max \left\{ x : (x, 1) \in S \right\} \\ &\stackrel{\text{U.B.}}{\leq} \mathbb{P}[b_0 < Q_0] + \mathbb{P}[b_1 > Q_1] \end{aligned}$$

By construction (of Q_0, Q_1), we have

$$\mathbb{P}[b_0 < Q_0] = D^m \left(\left\{ x : (x, y) \in S, x \notin (Q_0, Q_1) \right\} \right) = (1 - \varepsilon)^m$$

$$\Rightarrow \mathbb{P}[b_0 < Q_0] = (1 - \varepsilon)^m$$

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

Combined, we get: $\mathbb{P}[L_{D,f}(h_S) > \varepsilon] \leq \frac{2e^{-\varepsilon m}}{\delta} \Rightarrow m > \frac{\log\left(\frac{2}{\delta}\right)}{\varepsilon}$

No Free-Lunch Theorem

Learning task is defined via a distribution over $X \times Y$
 (in our case $X \times \{0,1\}$). Goal: Find a predictor

$$h: X \rightarrow Y = \{0,1\}$$

that has small risk, i.e.,

$$L_D(h) \text{ is small!}$$

Question: Is prior knowledge really necessary?

So far, prior knowledge come in the form of

1 realizability ($\exists h$ in some H , with $L_D(h)=0$)

2 assuming $\min_{h \in H} L_D(h)$ is small

(3 assuming something about D)

Thm: Let A be any learning algorithm for the task of binary classification (under 0-1 loss) over a domain X . Let m be any number smaller than $|X|/2$, representing the training set size, i.e., $|S|=m$. Then, there exists a distribution D over $X \times Y = \{0,1\}$, such

that 1 $\exists f: X \rightarrow \{0,1\}$ with $L_D(f)=0$!

2 with prob. of at least $\frac{1}{4}$ over the choice of $S \sim D^m$, we have $L_D(A(S)) \geq \frac{1}{8}$!

Wrt. [1], we could have, e.g. $H = \{f\}$, or H that is finite and contains f .

Proof: Let X be our domain and $C \subset X$, such that $|C| = 2^m$.

We see that there are 2^{2^m} functions from $C \rightarrow \{0,1\}$.

Let f_1, \dots, f_T denote these functions, $T = 2^{2^m}$.

Define:

$$D_i(\{(x,y)\}) = \begin{cases} \frac{1}{|C|}, & \text{if } y = f_i(x) \\ 0, & \text{else} \end{cases}$$

$\Rightarrow L_{D_i}(f_i) = 0$ by construction!

We show the following: For every algorithm A , receiving training data of size m from $X \times \{0,1\}$, it holds that

$$\max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(s))] \geq \frac{1}{4}$$

Consequence: For every algorithm A' , receiving m training samples from $X \times \{0,1\}$, there exists a function $f: X \rightarrow \{0,1\}$ and a distribution D over $X \times \{0,1\}$ such that $L_D(f) = 0$ and

$$\mathbb{E}_{S \sim D^m} [L_D(A'(s))] \geq \frac{1}{4}$$

Notation: $[T] = \{1, \dots, T\}$

Our training data, S , is $S = (x_1, \dots, x_m)$. Since $|C|=2^m$, there are $(2^m)^m$ possible training sequences from C .

Denote these sequences by S_1, \dots, S_k with $k = (2^m)^m$.

If S_j is labeled by f_i , we write S_j^i .

$$S_j^i = ((x_1, f_i(x_1)), \dots, (x_m, f_i(x_m)))$$

We can write

$$\mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(S))] = \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(S_j^i))$$

$$\Rightarrow \max_{i \in [T]} \mathbb{E}_{S \sim D_i^m} [L_{D_i}(A(S))] = \max_{i \in [T]} \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(S_j^i)) \quad (\star)$$

$$\begin{aligned} (\star) \quad \max_{i \in [T]} \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(S_j^i)) &\geq \frac{1}{T} \cdot \sum_{i=1}^T \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(S_j^i)) \\ &= \underbrace{\frac{1}{k} \cdot \sum_{j=1}^k}_{\text{avg. over } j=1, \dots, k} \underbrace{\frac{1}{T} \cdot \sum_{i=1}^T L_{D_i}(A(S_j^i))}_{\text{avg. over } i=1, \dots, T} \end{aligned}$$

$$\geq \min_{j \in [k]} \frac{1}{T} \cdot \sum_{i=1}^T L_{D_i}(A(S_j^i))$$

output is a hyp. h: $X \rightarrow \{0, 1\}$

Next step: Fix $j \in [k]$.

If $S = (x_1, \dots, x_m)$ and $|C| = 2^m$, there are $P \geq m$ remaining instances in C , denoted by v_1, \dots, v_P .

For any $h: C \rightarrow \{0, 1\}$, we have that

$$\begin{aligned} L_{D_i}(h) &= \frac{1}{2^m} \cdot \sum_{x \in C} \mathbb{1}_{h(x) \neq f_i(x)} \\ &\geq \frac{1}{2^m} \sum_{r=1}^P \mathbb{1}_{h(v_r) \neq f_i(v_r)} \\ &\geq \frac{1}{2^P} \sum_{r=1}^P \mathbb{1}_{h(v_r) \neq f_i(v_r)} \quad (\text{xx}) \end{aligned}$$

By using (xx), we get that

$$\begin{aligned} \frac{1}{T} \cdot \sum_{i=1}^T L_{D_i}(A(S_j^i)) &\geq \frac{1}{T} \cdot \sum_{i=1}^T \frac{1}{2^P} \cdot \sum_{r=1}^P \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)} \\ &= \frac{1}{2^P} \cdot \sum_{r=1}^P \frac{1}{T} \cdot \sum_{i=1}^T \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)} \\ &= \underbrace{\frac{1}{2} \cdot \frac{1}{P} \cdot \sum_{r=1}^P \frac{1}{T} \cdot \sum_{i=1}^T \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)}}_{\text{Avg. over } r=1, \dots, P} \\ (\text{xxx}) &\geq \frac{1}{2} \cdot \min_{r \in [P]} \frac{1}{T} \cdot \sum_{i=1}^T \mathbb{1}_{A(S_j^i)(v_r) \neq f_i(v_r)} \end{aligned}$$

Next step: Fix $r \in [p]$

$$[p] = \{1, \dots, p\}$$

C is of size $2m$.

$$\text{Soy } C = 4 = 2m \Rightarrow m = 2$$

	f_1, f_2, \dots	f_i
x_1	0 0 0 0 0 0 1 1 1 1 1 1 1 1	
x_2	0 0 0 0 1 1 1 1 0 0 0 1 1 1	
\dots		
$v_r = x_3$	0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1	
x_4	1 0 1 0 1 0 1 0 1 0 1 0 1 0 1	
	(brace)	

In general, we can partition all functions into $\frac{T}{2}$ disjoint pairs. For each pair $(f_i, f_{i'})$ we have that for every $c \in C$, $f_i(c) \neq f_{i'}(c) \Leftrightarrow c = v_r$.

It follows that $S_j^i = S_j^{i'}$ for each pair!

$$\Rightarrow \frac{1}{T} \sum_{j=1}^T A(S_j^i)(v_r) \neq f_i(v_r) + \frac{1}{T} \sum_{j=1}^T A(S_j^{i'})(v_r) \neq f_{i'}(v_r) - 1$$

One of them has to be 1, the other 0!

We get that

$$\frac{1}{T} \cdot \sum_{i=1}^T \frac{1}{T} \sum_{j=1}^T A(S_j^i)(v_r) \neq f_i(v_r) = \frac{1}{2}$$

With $(\times \times)$, we have

$$\frac{1}{T} \cdot \sum_{i=1}^T L_{D_i}(A(S_j^i)) \geq \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

In combination:

$$\max_{i \in [T]} \frac{1}{k} \cdot \sum_{j=1}^k L_{D_i}(A(s_j^i)) \geq \frac{1}{4}$$

$\underbrace{\mathbb{E}_{s \sim D_i^m} [L_{D_i}(A(s))]}$

\Rightarrow This shows that

$$\max_{i \in [T]} \mathbb{E}_{s \sim D_i^m} [L_{D_i}(A(s))] \geq \frac{1}{4}$$

and we know that this implies (first page) for any A' receiving m samples from $X \times \{0,1\}$ there exists a function $f: X \rightarrow \{0,1\}$ and distribution D over $X \times \{0,1\}$ such that $L_D(f) = 0$ and

$$\mathbb{E}_{s \sim D^m} [L_D(A'(s))] \geq \frac{1}{4}$$

(xxx)

Lemma (related to the MARKOV inequality). Let Z be a random variable that takes on values in $[0,1]$. Assume $\mathbb{E}[Z] = \mu$. Then, for any $\alpha \in (0,1)$

$$P[Z \geq 1-\alpha] \geq \frac{\mu - (1-\alpha)}{\alpha}$$

We needed to show that

$$\underset{s \sim D^m}{\mathbb{P}} \left[L_D(A(s)) \geq \frac{1}{8} \right] \geq \frac{1}{7} \quad \left(\begin{matrix} [12] \text{ in the} \\ \text{Thm.} \end{matrix} \right)$$

We write (xxxx)

$$\underset{s \sim D^m}{\mathbb{P}} \left[L_D(A(s)) \geq 1 - \underbrace{\frac{7}{8}}_{\text{will be our } \alpha \in (0,1)} \right] \geq \frac{\mathbb{E}[L_D(A(s))] - (1 - \frac{7}{8})}{7/8}$$
$$\geq \frac{1/8}{7/8}$$
$$= \frac{1}{7}$$

No prior knowledge!

Corollary: Let X be an infinite domain and let H be the set of all functions from $X \rightarrow \{0,1\}$. Then H is not PAC learnable.



VC - Dimension

(Vapnik & Chervonenkis, 1970)

We have seen that without any restrictions on the hyp. class, an adversary can always construct a distribution on which the learning alg. will **perform poorly**, while another alg. succeeds on the same distribution.

Why: We had the "freedom" of choosing a target function (i.e., one of the f_i) from **all** possible functions from $C \rightarrow \{0, 1\}$.

Let's say we have an hyp. class H now. For PAC learnability, the adversary is now limited to the construction of distributions for which $\exists h^* \in H$ with **zero risk!**

As a starting point, given a finite $C \subset \mathcal{X}$, we consider how H behaves on C (note that in the No-Free-Lunch theorem, the distributions were concentrated on C).

Def. Restricting
 H to C

Let H be a hyp. class of functions from $X \rightarrow \{0,1\}$ and $C = \{c_1, \dots, c_m\} \subset X$. The **restriction of H to C** is

$$H_C = \{(h(c_1), \dots, h(c_m)) : h \in H\}$$

H_C is a set of functions, each represented by a vector $(0, 1, 0, 0, \dots, 1)$ of length $|C|$. In other words, these are all labelings of C that are possible using the $h \in H$.

Def. Shattering

For a finite set $C \subset X$, if H_C contains **all** functions from $C \rightarrow \{0,1\}$, we say **H shatters C** .

Ex.: Say $C = \{c_1, c_2\}$. If $H_C = \{(0,0), (0,1), (1,0), (1,1)\}$ then H shatters C , as there are 4 possible labelings of C .

More concrete example: Threshold functions over \mathbb{R} .

Visually: $\text{Label } 1$ $\text{Label } 0$ \mathbb{R}

$\hookrightarrow \text{Threshold } \in \mathbb{R}$

[1] Start with $C = \{c_1\}$

$$|C|=1$$

If $q = c_1 + 1$, then $h_q(c_1) = 1$

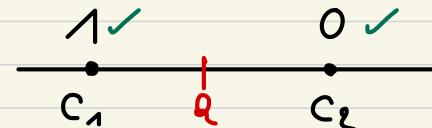
If $q = c_1 - 1$, then $h_q(c_1) = 0$

\Rightarrow as we can choose q freely, we see that we can realize all possible labelings of $C = \{c_1\}$ and thus H shatters C .

[2] We go on with $C = \{c_1, c_2\}$, $|C|=2$, and let $c_1 \leq c_2$ without loss of generality (w.l.o.g.).

Visual argument:

$\Rightarrow H$ does NOT shatter C ! As c_1, c_2 were arbitrary, H does not shatter any set of size 2!



NOT POSSIBLE!

$\hookrightarrow H_C$ is $\{(1,0), (1,1), (0,0)\} \Rightarrow |H_C|=3 \neq 2^{|C}|=4!$

Corollary (of
the No-Free -
Lunch thm.)
(*)

Let H be a hyp. class of functions from $X \rightarrow \{0,1\}$.
let m be a training set size. Assume there exists a
set C of size $2m$ that is shattered by H . Then, for any
learning alg. A , there exist a distribution D over $X \times \{0,1\}$
and a predictor $h \in H$ such that $L_D(h) = 0$, but with
probability of at least $1/4$ over the choice $S \sim D^m$, we
have $L_D(A(S)) \geq 1/8$.

Remark:

Note that the fact that H shatters C of size $2m$
means that H can realize all labelings of the $2m$ points
in C . Hence, we are exactly in the setting of the No-
Free-Lunch theorem, as our alg. only receives samples
 S of size m !

Consequently, if H is such that it shatters a large
set C , then we cannot learn by just seeing half of
the samples.

Def. VC-dimension.

The VC-dimension of a hyp. class H , denoted by $VC(H)$

is the maximum size of a set $C \subset X$ that can be shattered by H .

Note: If H shatters sets of arbitrary size, we say that H has infinite VC-dimension ($VC(H) = \infty$)

Corollary:

Let H be a hyp. class with $VC(H) = \infty$. Then H is not learnable.

Proof: No matter what training set size m we have, there is a set of size $2m$ that is shattered. Hence, by our corollary (*) from the previous page, the claim follows. □

(we will see that the converse also holds, i.e., finite VC-dimension implies PAC learnability).

Ex. Threshold functions on \mathbb{R} revisited: we saw that H shatters $C = \{c_1\} \Rightarrow VC(H) \geq 1$. Also, H does NOT shatter sets $C = \{c_1, c_2\}$ of size 2 $\Rightarrow VC(H) < 2$. Thus $VC(H) = 1$.

General strategy
to show the VC
dim. of H .
 $(VC(H) = d)$

$$d < \infty$$

- [1] There exists a set C of size d that is shattered.
- [2] Every set of size $d+1$ is not shattered by H .

What about the VC-dim. of finite hyp. classes?

If H is finite, then for any set C , we have

$$|H_C| \leq |H|$$

E.g., if we only have 2 hypotheses and a set C , we can, at most, create 2 different labelings of C . If H shatters C , we know that $|H_C| = 2^{|C|}$ and $|C|$ is the max. size that is shattered, i.e., the VC-dim. We get

$$|H| \geq |H_C| = 2^{|C|} = 2^{VC(H)}$$

$$\Rightarrow \log_2(|H|) \geq VC(H)$$

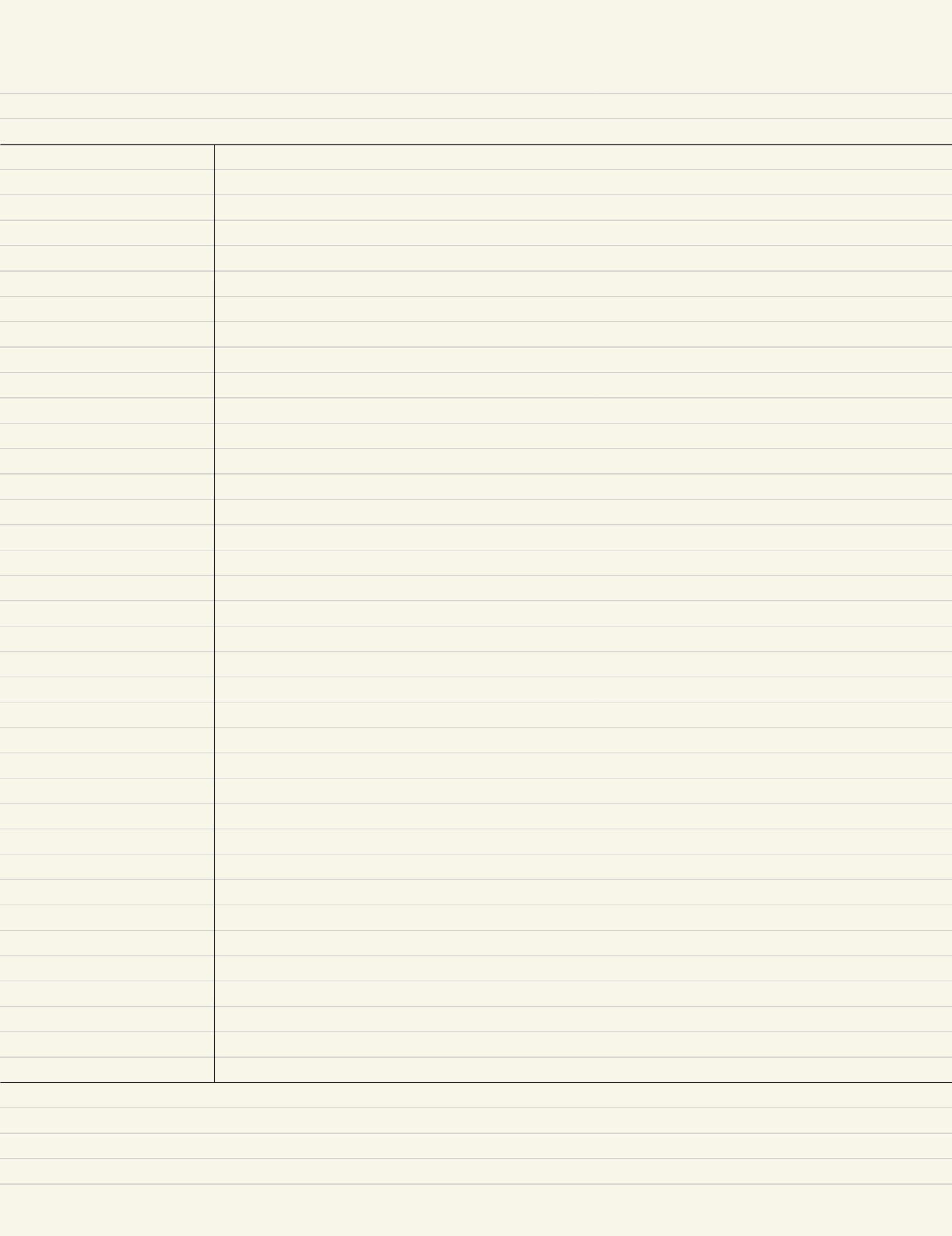
$$VC(H) \leq \log_2(|H|)$$

Note that this upper-bound can be quite loose!

E.g., k threshold functions $\rightarrow |H| = k$ but $VC(H) = 1$!

Finally, we remark that $VC(H) = d$ means that there exists a set of size d that is shattered. This DOES NOT imply that ALL sets of size d or SMALLER are fully shattered (again, there just exists one of size d that is!).

[I will cover some examples next, before we move on
to connect all these things together!]



Plan: Proof of $\boxed{16} \rightarrow \boxed{1}$

Def. (Growth)
function

Let H be a hypothesis class. Then, the growth function of H , denoted by $\mathbb{J}_H : \mathbb{N} \rightarrow \mathbb{N}$, is defined as

$$\mathbb{J}_H(m) = \max_{C \subseteq X : |C|=m} |H_C|$$

This means, $\mathbb{J}_H(m)$ counts the number of different functions from a set C (of size m) to $\{0, 1\}$.

SAUER'S LEMMA:

Let H be a hyp. class with $VC(H) = d < \infty$. Then, for all m

$$\mathbb{J}_H(m) \leq \sum_{i=0}^d \binom{m}{i}$$

\Rightarrow exponential growth with m !

In particular, if $m > d + 1$, we have

$$\mathbb{J}_H(m) \leq \left(\frac{e \cdot m}{d}\right)^d$$

\Rightarrow polynomial growth with m !

Remark: if $VC(H) = d$, then for $m \leq d$, we have $\mathbb{J}_H(m) = 2^m$.

Why? Because we can shatter a set of size $d \Rightarrow$ all labelings of the points on that set are realizable via functions in H).

Theorem :

Let H be a hyp. class and γ_H its growth function. Then, for every distribution D and every $\delta \in (0, 1)$, with probability of at least $1 - \delta$ (over the choice of S , $S \sim D^m$), we have for a loss function with range $[0, c]$

$$\forall h \in H: |L_D(h) - L_S(h)| \leq c \sqrt{\frac{8 \cdot \log(\gamma_H(2m) \cdot \delta^{-1})}{m}}$$

Proof: Recall that

$$|a - b| \geq |a| - |b| \quad (\text{triangle inequality})$$

and

$$L_D(h) = \mathbb{E}_{S' \sim D^m} [L_{S'}(h)] \quad (\text{def})$$

$$|(L_D(h) - L_S(h)) - (L_D(h) - L_{S'}(h))| \stackrel{(1)}{\geq} |L_D(h) - L_S(h)| - |L_D(h) - L_{S'}(h)|$$

This means that

$$|L_{S'}(h) - L_S(h)| \geq \underbrace{|L_D(h) - L_S(h)|}_{(a)} - \underbrace{|L_D(h) - L_{S'}(h)|}_{(b)}.$$

If $(a) > \varepsilon$ and $(b) < \frac{\varepsilon}{2}$, we get

$$|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}$$

Write (a) and (b) in terms of indicator functions. In fact,
 $|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}$ only holds if (a) and (b) are satisfied.

Hence,

$$\mathbb{P} |L_D(h) - L_S(h)| > \varepsilon \cdot \underbrace{\mathbb{P} |L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2}}_{\boxed{T1}} \leq \underbrace{\mathbb{P} |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}}_{\boxed{T2}}$$

$\boxed{T1}$ Take $\mathbb{E}_{S'}[\cdot]$:

$$\begin{aligned} \mathbb{E}_{S'} \left[\mathbb{P} |L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2} \right] &= \mathbb{P} \left[|L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2} \right] \\ (\text{Hoeffding}) &\geq 1 - 2e^{-2m(\frac{\varepsilon}{2})^2 / (S - Q)^2} \rightarrow c^2 \\ &= 1 - 2e^{-m\varepsilon^2 / 2c^2} \end{aligned}$$

$\boxed{T2}$ Take $\mathbb{E}_{S'}[\cdot]$:

$$\begin{aligned} \mathbb{E}_{S'} \left[\mathbb{P} |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] &= \mathbb{P}_{S'} \left[|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] \\ &\leq \mathbb{P}_{S'} \left[\exists h \in H : |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] \end{aligned}$$

Assumption:
(we will check this later)

$$m \geq 4c^2 \cdot \epsilon^{-2} \cdot \log(2) \quad (\text{xxx})$$

Say we are at equality (=). Then,

$$\begin{aligned} \Pr_{S'} \left[|L_D(h) - L_{S'}(h)| < \frac{\epsilon}{2} \right] &\geq 1 - 2e^{-m\epsilon^2/2c^2} \quad (\text{now, set } c \\ &= \dots \text{ (xxx)}) \\ &= 1 - 2e^{-(4c^2\epsilon^{-2}\log(2))\epsilon^2/2c^2} \\ &= 1 - 2e^{-2\cdot\log(2)} \\ &= 1 - 2e^{-\log(2^2)} \\ &= 1 - 2e^{-\log(4)} \\ &= 1 - \frac{2}{4} = 1 - \frac{2}{4} = \frac{1}{2} \end{aligned}$$

In combination with $\mathbb{1}_{|L_D(h) - L_S(h)| > \epsilon}$, we get

$$(\text{xxx}) \quad \mathbb{1}_{|L_D(h) - L_S(h)| > \epsilon} \leq 2 \cdot \Pr_{S'} \left[\exists h \in H: |L_{S'}(h) - L_S(h)| > \frac{\epsilon}{2} \right]$$

(under the assumption on $m \geq \dots$ from above).

In fact (xxx) holds for **all $h \in H$!** Hence, we can write

$$\mathbb{1}_{\exists h \in H: |L_D(h) - L_S(h)| > \epsilon} \leq 2 \cdot \Pr_{S'} \left[\exists h \in H: |L_{S'}(h) - L_S(h)| > \frac{\epsilon}{2} \right]$$

Let's take $\mathbb{E}_S [\cdot]$ now!

we get

$$\mathbb{P}_{S} \left[\exists h \in H : |L_D(h) - L_S(h)| > \frac{\epsilon}{2} \right] \leq 2 \cdot \mathbb{P}_{S, S'} \left[\exists h \in H : |L_{S'}(h) - L_S(h)| > \frac{\epsilon}{2} \right]$$

elements in S or S'

we only have empirical risks!

TRICK: Note that \hat{z}_i, z_i are drawn IID from D . Hence, we can easily switch them. What effect do we get?

$$(L(h, \hat{z}_i) - L(h, z_i)) \xrightarrow{\text{(upon switch)}} - (L(h, z_i) - L(h, \hat{z}_i))$$

$L(h, \hat{z}_i)$	$L(h, z_i)$	$(L(h, \hat{z}_i) - L(h, z_i))$
0	0	0
0	1	-1
1	0	1
1	1	0

In $|L_{S'}(h) - L_S(h)|$, we actually have

$$\left| \frac{1}{m} \cdot \sum_{i=1}^m (L(h, \hat{z}_i) - L(h, z_i)) \right|$$

Let's use this insight!

$$\begin{aligned}
 & \mathbb{P}_{S, S'} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] \leq 2 \cdot \mathbb{P}_{S, S'} \left[\exists h \in H : |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] \\
 &= 2 \cdot \mathbb{P}_{S, S'} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m (\ell(h, z_i') - \ell(h, z_i)) \right| > \frac{\varepsilon}{2} \right] \\
 &= 2 \cdot \mathbb{E}_{S, S'} \left[\# \left\{ \exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m (\ell(h, z_i') - \ell(h, z_i)) \right| > \frac{\varepsilon}{2} \right\} \right] \\
 &= 2 \cdot \mathbb{E}_{S, S'} \left[\# \left\{ \# \left\{ \exists h \in H : \frac{1}{m} \cdot \sum_{i=1}^m \Theta_i (\ell(h, z_i') - \ell(h, z_i)) \right| > \frac{\varepsilon}{2} \right\} \right]
 \end{aligned}$$

We can do this if Θ_i is drawn from a uniform distribution on $\{+1, -1\}$. So $\Theta_i \sim U(\{-1, +1\})$.

↳ uniform distr.

Note: The Θ_i 's are called RADEMACHER variables!

We have now that

$$\dots \leq 2 \cdot \mathbb{E}_{S, S'} \left[\# \left\{ \exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m \Theta_i (\ell(h, z_i') - \ell(h, z_i)) \right| > \frac{\varepsilon}{2} \right\} \right]$$

As $\sum_{i=1}^m \dots$ goes over $1, \dots, m$ and $z_i', z_i \sim D$, we see that a $h \in H$ is only evaluated on at most $2m$ points! So, let's combine S and S' into a set C , $|C| \leq 2m$ and restrict H to C .

Since C is finite, we have that H_C is also finite!

Let's look at the $\underset{\sigma}{\mathbb{P}}[\dots]$ term. by union bound

$$\underset{\sigma}{\mathbb{P}} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m \sigma_i \cdot (\ell(h, z_i) - \ell(h, t_i)) \right| > \frac{\varepsilon}{2} \right] \leq$$

$$\sum_{h \in H_C} \underset{\sigma}{\mathbb{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{1}{m} \cdot (\ell(h, z_i) - \ell(h, t_i)) \right| > \frac{\varepsilon}{2} \right] =$$

$$\sum_{h \in H_C} \underset{\sigma}{\mathbb{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{1}{m} \cdot k_i \right| > \frac{\varepsilon}{2} \right]$$

$$\underset{\sigma}{\mathbb{E}} \left[\mathbb{1}_{\left| \sum \sigma_i \cdot \frac{1}{m} \cdot k_i \right| > \frac{\varepsilon}{2}} \right]$$

As all σ_i are drawn IID from $\mathcal{U}\{\pm 1\}$, we know that the expectation is 0 !

$$\rightarrow \sum_{h \in H_C} \underset{\sigma}{\mathbb{E}} \left[\mathbb{1}_{\left| \sum \sigma_i \cdot \frac{1}{m} \cdot k_i \right| > \frac{\varepsilon}{2}} \right] = \sum_{h \in H_C} \underset{\sigma}{\mathbb{E}} \left[\mathbb{1}_{\left| \sum \sigma_i \cdot \frac{1}{m} \cdot k_i - 0 \right| > \frac{\varepsilon}{2}} \right]$$

$$= \sum_{h \in H_C} \underset{\sigma}{\mathbb{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{1}{m} \cdot k_i - 0 \right| > \frac{\varepsilon}{2} \right] \quad (\text{use Hoeffding inequality})$$

Side result: $\underset{\sigma}{\mathbb{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{1}{m} \cdot k_i - 0 \right| > \frac{\varepsilon}{2} \right] \leq 2 \cdot e^{-2m \left(\frac{\varepsilon}{2} \right)^2 / (2c)^2}$

$$= 2 \cdot e^{-2m \varepsilon^2 / 8c^2}$$

$$= 2 \cdot e^{-m \varepsilon^2 / 8c^2}$$

Combined, we get:

$$\underset{\sigma}{\mathbb{P}} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_i \sigma_i \cdot (\ell(h, z_i) - \ell(h, t_i)) \right| > \frac{\varepsilon}{2} \right] \leq 2 \cdot e^{-m \varepsilon^2 / 8c^2} \cdot |H_C|$$

Because we have $\sum_{h \in H_C} \dots$

Finally :)

$$\underset{S}{\mathbb{P}} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] \leq 2 \cdot \underset{S, S'}{\mathbb{E}} \left[|H_C| \cdot 2 \cdot e^{-m \cdot \varepsilon^2 / 8c^2} \right]$$

$$= 4 \cdot \underset{S, S'}{\mathbb{E}} \left[|H_C| \cdot e^{-m \cdot \varepsilon^2 / 8c^2} \right]$$

Now, we replace $|H_C|$ by $\gamma_H(2m)$ \Rightarrow we no longer have a dependency on S or S' !

Overall, we get

$$\underset{S}{\mathbb{P}} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] \leq \underbrace{4 \cdot \gamma_H(2m)}_{\leq \delta} \cdot e^{-m \cdot \varepsilon^2 / 8c^2}$$

$$4 \gamma_H(2m) e^{-m \varepsilon^2 / 8c^2} < \delta$$

$$e^{-m \varepsilon^2 / 8c^2} < \frac{\delta}{4 \gamma_H(2m)} \quad | \log$$

$$-m \varepsilon^2 / 8c^2 < \log \left(\frac{\delta}{4 \gamma_H(2m)} \right)$$

$$m \varepsilon^2 > 8c^2 \cdot \log \left(\frac{4 \gamma_H(2m)}{\delta} \right)$$

$$\varepsilon > \sqrt{\frac{8 \cdot \log(4 \gamma_H(2m) / \delta)}{m}}$$

This is the ε -term in our theorem!

Let's check the assumption!

$$\text{Recall: } m \geq 4c^2 \epsilon^{-2} \log(2) \quad (\text{Set in } \varepsilon)$$

$$\Rightarrow m \geq 4c^2 \cdot \frac{1}{\epsilon \cdot \frac{8 \log(4\gamma_H(\ell_m)/\delta)}{m}} \cdot \log(2)$$

$$8 \log\left(\frac{4\gamma_H(\ell_m)/\delta}{\epsilon}\right) \geq 4 \cdot \log(2)$$

$$\Rightarrow \log(\dots) \geq \frac{1}{2} \cdot \log(2) \quad | e^{..}$$

$$\Rightarrow 4\gamma_H(\ell_m)/\delta \geq \underbrace{c^{\frac{1}{2} \cdot \log(2)}}_{\sqrt{2}}$$

$$\Rightarrow \frac{1}{\delta} \geq \frac{\sqrt{2}}{4\gamma_H(\ell_m)}$$

$$\Rightarrow \delta \leq 4 \cdot \gamma_H(\ell_m) \cdot \frac{1}{\sqrt{2}}$$

$$\Rightarrow \delta \leq 2\sqrt{2} \cdot \gamma_H(\ell_m)$$

δ is in $(0, 1) \rightarrow$ so, the last inequality is always satisfied!

We have shown that $P_S \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] < \delta$ (for ε as before). This is equivalent to

$$P_S \left[\forall h \in H : |L_D(h) - L_S(h)| \leq \varepsilon \right] \geq 1 - \delta$$

qed

16] → 17] Tricky!

Our proof will be different to the book, but show the same claim.

Before, we introduce a very important result in learning theory.

Def. growth function

let H be a hyp. class. Then the growth function of H , denoted by $\gamma_H : \mathbb{N} \rightarrow \mathbb{N}$, is defined as

$$\gamma_H(m) = \underbrace{\max_{C \subset X: |C|=m} |H_C|}_{\# \text{ of different functions from a set } C \text{ of size } m \text{ to } \{0,1\}}$$

Observation. If $VC(H) = d$, then for $m \leq d$, we have
 $\gamma_H(m) = 2^m$ (as we can shatter a set of size $d \Rightarrow$ all labelings of d points are realizable w.r.t. functions in H)

Let H be a hyp. class (as before) with $VC(H) = d < \infty$.

Then, for all m $\gamma_H(m) \leq \sum_{i=0}^d \binom{m}{i}$ (Exponential growth with m .)

In particular, if $m > d+1$, we have

$$\gamma_H(m) \leq \left(\frac{em}{d}\right)^d \quad (\text{Polynomial growth with } m.)$$

Sauer's lemma

while, we will NOT proof Sauer's lemma, we proof $\boxed{G} \rightarrow \boxed{1}$ now.

Thm.

Let H be a hyp. class and γ_H it's growth function. Then, for every distribution D and every $\delta \in (0, 1)$, with prob. of at least $1-\delta$ over the choice of $S \sim D^m$, we have for a loss function in the range $[0, c]$

$$\forall h \in H: |L_D(h) - L_S(h)| \leq c \cdot \sqrt{\frac{8 \cdot \log(\gamma_H(2m))^2}{m}} \delta$$

Proof (sorry, it's quite long):

Let's start with the triangle inequality for absolute values:

$$|a - b| \geq |a| - |b| \quad (\times)$$

Also, recall that

$$L_D(h) = \mathbb{E}_{S^1 \sim D^m} [L_{\delta^1}(h)] \quad (\times \times)$$

This allows us to replace $L_D(h)$ in the theorem by $(\times \times)$. Assume we have $S, S' \sim D^m$, i.e., two random samples of size m . Using (\times) , we have

$$\begin{aligned} |(L_D(h) - L_S(h)) - (L_D(h) - L_{S'}(h))| &\geq |L_D(h) - L_S(h)| - \\ &|L_D(h) - L_{S'}(h)| \end{aligned}$$

Re-arranging terms gives

$$|L_{S'}(h) - L_S(h)| \geq \underbrace{|L_D(h) - L_S(h)|}_{(a)} - \underbrace{|L_D(h) - L_{S'}(h)|}_{(b)}$$

If (a) > ε and (b) < $\frac{\varepsilon}{2}$, we get

$$|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}$$

We could try to write these conditions on (a) and (b) in terms of **indicator functions** and formulate an inequality:

In fact $|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}$ only holds if both (a) and (b) are satisfied. Hence

$$\underbrace{\mathbb{E} [\mathbb{1}_{|L_D(h) - L_S(h)| > \varepsilon}]}_{\boxed{T1}} \cdot \underbrace{\mathbb{E} [\mathbb{1}_{|L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2}}]}_{\boxed{T2}} \leq \mathbb{E} [|L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2}]$$

$\boxed{T1}$ if we take $\mathbb{E}_{S'}[\cdot]$, we get

$$\begin{aligned} \mathbb{E}_{S'} \left[\mathbb{1}_{|L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2}} \right] &= \mathbb{P}_{S'} \left[|L_D(h) - L_{S'}(h)| < \frac{\varepsilon}{2} \right] \\ (\text{by Hoeffding}) &\geq 1 - 2e^{-2m(\frac{\varepsilon}{2})^2/(b-a)^2} \\ &= 1 - 2e^{-m\varepsilon^2/2c^2} \end{aligned}$$

The other term including S^1 is $\boxed{\mathbb{P}_{S^1}}$. A priori, $\mathbb{E}_{S^1}[\cdot]$ gives

$$\mathbb{E}_{S^1} \left[\frac{1}{2} |L_{S^1}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] = \mathbb{P}_{S^1} \left[|L_{S^1}(h) - L_S(h)| > \frac{\varepsilon}{2} \right]$$

$$\leq \mathbb{P}_{S^1} \left[\exists h \in H : |L_{S^1}(h) - L_S(h)| > \frac{\varepsilon}{2} \right]$$

Assumption: $m \geq 4c^2 \bar{\varepsilon}^2 \log(2)$

Say, we are at equality, i.e., $m = \dots$, then

$$\begin{aligned} \mathbb{P}_{S^1} \left[|L_D(h) - L_{S^1}(h)| < \frac{\varepsilon}{2} \right] &\geq 1 - 2e^{-\frac{m\varepsilon^2}{2c^2}} \\ &\quad - (4c^2 \bar{\varepsilon}^2 \log(2)) \frac{\varepsilon^2}{2c^2} \\ &= 1 - 2e^{-\frac{m\varepsilon^2}{2c^2}} \\ &= 1 - 2e^{-\frac{2 \log(2)}{c^2}} \\ &= 1 - 2e^{-\frac{\log(2^2)}{c^2}} \\ &= 1 - 2e^{-\frac{1}{c^2}} = 1 - \frac{1}{2} = \frac{1}{2} \end{aligned}$$

In combination with $\frac{1}{2} |L_D(h) - L_S(h)| > \varepsilon$ (which is not affected by $\mathbb{E}_{S^1}[\cdot]$), we get

$$\frac{1}{2} |L_D(h) - L_S(h)| > \varepsilon \leq 2 \cdot \mathbb{P}_{S^1} \left[\exists h \in H : |L_{S^1}(h) - L_S(h)| > \frac{\varepsilon}{2} \right]$$

(under the assumption on $m \geq \dots$)

Note that the previous inequality holds for all h . Hence, we can write

$$\frac{1}{n} \sum_{h \in H} |L_D(h) - L_S(h)| > \varepsilon \leq 2 \cdot \overline{P}_{S, S'} \left[\sum_{h \in H} |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right]$$

Let's take $\overline{E}[\cdot]$ now:

$$\overline{P}_{S'} \left[\sum_{h \in H} |L_D(h) - L_S(h)| > \varepsilon \right] \leq 2 \cdot \overline{P}_{S, S'} \left[\sum_{h \in H} |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right]$$

Both are empirical sums

Now comes a TRICK:

Note that all z_i, z'_i are drawn IID from D . Hence, we can easily switch them. So, what changes?

$$(\ell(h, z'_i) - \ell(h, z_i)) \rightarrow -(\ell(h, z'_i) - \ell(h, z_i))$$

why?

$\ell(h, z'_i)$	$\ell(h, z_i)$	$(\ell(h, z'_i) - \ell(h, z_i))$
0	0	0
0	1	-1
0	0	0
1	0	1

\Rightarrow switching z'_i and $z_i \Rightarrow$ sign change

In $|L_{S'}(h) - L_S(h)|$, we actually lose

$$\left| \frac{1}{m} \cdot \sum_{i=1}^m (l(h, z_i') - l(h, z_i)) \right|$$

if we would multiply by -1 at a specific i , we would effectively switch z_i' and z_i !

Let's use this insight:

$$\begin{aligned} \underset{S}{\mathbb{P}} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] &\leq 2 \cdot \underset{S, S'}{\mathbb{P}} \left[\exists h \in H : |L_{S'}(h) - L_S(h)| > \frac{\varepsilon}{2} \right] \\ &= 2 \cdot \underset{S, S'}{\mathbb{P}} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m (l(h, z_i') - l(h, z_i)) \right| > \frac{\varepsilon}{2} \right] \\ &= 2 \cdot \underset{S, S'}{\mathbb{E}} \left[\mathbb{1}_{\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m (l(h, z_i') - l(h, z_i)) \right| > \frac{\varepsilon}{2}} \right] \\ &= 2 \cdot \underset{S, S'}{\mathbb{E}} \left[\underset{\sigma}{\mathbb{E}} \left[\mathbb{1}_{\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m (\underbrace{l(h, z_i') - l(h, z_i)}_{\sigma_i}) \right| > \frac{\varepsilon}{2}} \right] \right] \end{aligned}$$

We can do this, if σ_i is drawn from a uniform distribution on $\{-1, +1\}$. The σ_i are so called RADEMACHER variables.

$$\Rightarrow 2 \cdot \underset{S, S'}{\mathbb{E}} \left[\underset{\sigma}{\mathbb{E}} \left[\mathbb{P} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m \sigma_i (l(h, z_i') - l(h, z_i)) \right| > \frac{\varepsilon}{2} \right] \right] \right]$$

Does not look easier, right? :)

BUT, as $\sum_{i=1}^m$ goes over $1, \dots, m$ and $z_i, z'_i \sim D$, we know that a $h \in H$ is only evaluated on at most $2m$ points!

In fact, we combine S and S' into C , $|C| \leq 2m$ and restrict H to C , i.e., H_C . Since C is finite, so is H_C . The question is how fast does H_C grow?

As always, we bound the \exists statement with the union bound (u.b.) and get (for the $\underset{\sigma}{\text{P}}[\cdot]$ term)

$$\begin{aligned} & \underset{\sigma}{\text{P}} \left[\exists h \in H : \frac{1}{m} \cdot \left| \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z'_i)) \right| > \frac{\epsilon}{2} \right] \stackrel{\text{u.b.}}{\leq} \\ & \sum_{h \in H_C} \underset{\sigma}{\text{P}} \left[\frac{1}{m} \cdot \left| \sum_{i=1}^m \sigma_i (\ell(h, z_i) - \ell(h, z'_i)) \right| > \frac{\epsilon}{2} \right] \\ & = \sum_{h \in H_C} \underset{\sigma}{\text{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{1}{m} \cdot (\ell(h, z_i) - \ell(h, z'_i)) \right| > \frac{\epsilon}{2} \right] \\ & \quad - \underbrace{\sum_{h \in H_C} \underset{\sigma}{\text{P}} \left[\left| \sum_{i=1}^m \sigma_i \cdot \frac{k_i}{m} \right| > \frac{\epsilon}{2} \right]}_{\mathbb{E}_{\sigma} \left[\mathbb{1} \left| \sum_{i=1}^m \sigma_i \cdot \frac{k_i}{m} \right| > \frac{\epsilon}{2} \right]} \end{aligned}$$

We know that the σ_i are drawn i.i.d from a uniform dist. on $\{-1, +1\} \Rightarrow$ expectation is 0!

We could (without changing anything) write

$$\sum_{h \in H_c} \mathbb{E}_\sigma \left[\mathbb{1}_{\left| \sum_{i=1}^m (\sigma_i \cdot \frac{k_i}{m} - 0) \right| > \frac{\epsilon}{2}} \right], \text{ or}$$

$$\sum_{h \in H_c} \mathbb{P}_\sigma \left[\left| \sum_{i=1}^m \frac{k_i}{m} \cdot \sigma_i - 0 \right| > \frac{\epsilon}{2} \right]$$

This seems to be a classic case for Hoeffding again, and we see why we introduced the RADIXCHARTER variables :)

Note that $k_i := (\ell(h, z'_i) - \ell(h, z_i))$ and we know the loss is in $[0, c]$ \Rightarrow the difference is in $[-c, c]$.

$$\begin{aligned} \Rightarrow \mathbb{P}_\sigma \left[\left| \sum_{i=1}^m \frac{k_i}{m} \sigma_i - 0 \right| > \frac{\epsilon}{2} \right] &= \mathbb{P}_\sigma \left[\left| \frac{1}{m} \sum_{i=1}^m k_i - 0 \right| > \frac{\epsilon}{2} \right] \\ &\leq 2 \cdot e^{-2m \cdot \left(\frac{\epsilon}{2}\right)^2 / (2c)^2} \\ &= 2 \cdot e^{-2m \epsilon^2 / 16c^2} = 2 \cdot e^{-m \epsilon^2 / (8c^2)} \end{aligned}$$

Combined with $\sum_{h \in H_c}$, we get

$$\mathbb{P}_\sigma \left[\exists h \in H_c : \frac{1}{m} \cdot \left| \sum_i (\ell(h, z'_i) - \ell(h, z_i)) \cdot \sigma_i \right| > \frac{\epsilon}{2} \right] \leq |H_c| \cdot 2e^{-m \epsilon^2 / 8c^2}$$

Finally :), we are ready to conclude that

$$\underset{S}{P} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] \leq 2 \cdot \underset{\delta, \delta'}{\mathbb{E}} \left[|H_\delta| \cdot 2 e^{-\frac{m\varepsilon^2}{8c^2}} \right]$$

$$= 4 \underset{\delta, \delta'}{\mathbb{E}} \left[|H_\delta| \cdot 2 e^{-\frac{m\varepsilon^2}{8c^2}} \right]$$

We are now ready to replace $|H_\delta|$ by $\gamma_H(2m)$, remembering that $|C| \leq 2m$. We get (as $\gamma_H(2m)$ is independent of δ, δ') :

$$\underset{S}{P} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] \leq \underbrace{4 \gamma_H(2m) \cdot e^{-\frac{m\varepsilon^2}{8c^2}}}_{< \delta}$$

$$4 \gamma_H(2m) \cdot e^{-\frac{m\varepsilon^2}{8c^2}} < \delta$$

$$e^{-\frac{m\varepsilon^2}{8c^2}} < \frac{\delta}{4 \gamma_H(2m)} \quad | \log.$$

$$-\frac{m\varepsilon^2}{8c^2} < \log \left(\frac{\delta}{4 \gamma_H(2m)} \right)$$

$$m\varepsilon^2 > 8c^2 \cdot \log \left(\frac{4 \gamma_H(2m)}{\delta} \right)$$

$$\varepsilon > c \sqrt{\frac{8 \cdot \log \left(\frac{4 \gamma_H(2m)}{\delta} \right)}{m}}$$

This is our Σ from the theorem's claim!

It remains to check if our assumption on $m \geq \dots$ was valid. In particular, this was

$$m \geq 4C^2 \bar{\epsilon}^2 \log(2)$$

$$m \geq 4\cancel{C^2} \cdot \frac{1}{\cancel{C} \cdot \frac{8 \ln(4J_H(2m)/\delta)}{m}} \cdot \log(2)$$

$$m \geq \frac{4 \cdot \cancel{m} \log(2)}{8 \log(4J_H(2m)/\delta)}$$

$$8 \log(4J_H(2m)/\delta) \geq 4 \cdot \log(2)$$

$$\log(\cdot) \geq \frac{1}{2} \cdot \log(2)$$

$$4J_H(2m)/\delta \geq e^{\frac{1}{2} \log(2)} = T_2$$

$$\frac{1}{\delta} \geq \frac{T_2}{4J_H(2m)}$$

$$\frac{1}{\delta} \leq \frac{4J_H(2m)}{\sqrt{T_2}}$$

$$\delta \leq 2\sqrt{T_2} \cdot J_H(2m)$$

As $\delta \in (0, 1)$, this last inequality is ALWAYS satisfied \Rightarrow OK!

Overall, we have shown that

$$\underset{s}{\mathbb{P}} \left[\exists h \in H : |L_D(h) - L_S(h)| > \varepsilon \right] < \delta \quad \text{→ we made this explicit}$$

which is equivalent to

$$\underset{s}{\mathbb{P}} \left[\forall h \in H : |L_D(h) - L_S(h)| \leq \varepsilon \right] \geq 1 - \delta$$

or, in other words "that with prob. of at least $1 - \delta$

$$\forall h \in H : |L_D(h) - L_S(h)| \leq c \cdot \sqrt{\frac{8 \cdot \log(4 \cdot |H| \cdot m) / \delta}{m}}$$

□

In the book, it is shown that the supremum over H is bounded, but since, in our case, we have $\forall h \in H$, this holds for the supremum as well!

! This establishes the last equivalence $\boxed{S} \rightarrow \boxed{H}$ in the fundamental theorem of PAC learning!

LINEAR PREDICTORS

Def.

We define the class of affine functions

$$L_d = \{ h_{w,b} : w \in \mathbb{R}^d, b \in \mathbb{R} \}$$

with

$$h_{w,b}(x) = \langle w, x \rangle + b = \left(\sum_{i=1}^d w_i x_i \right) + b \quad (x \in \mathbb{R}^d)$$

There are different hypothesis classes of linear predictors and each class is a composition of a function

$$\phi : \mathbb{R} \rightarrow Y$$

with L_d . Example: $\phi = \text{sign}$. In that case, i.e., $\phi = \text{sign}$, we get the hypothesis class of HALFSPACES.

Notation:

We will also write L_d as $\{ x \mapsto \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R} \}$

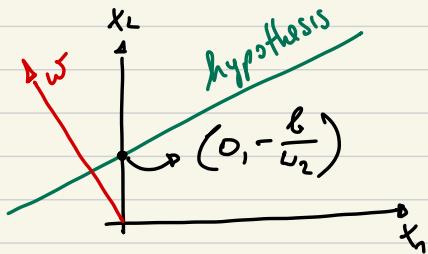
HALFSPACES

Let $X = \mathbb{R}^d$ and $Y = \{+1, -1\}$. Then

$$HS_d = \text{sign} \circ L_d = \{ x \mapsto \text{sign}(h_{w,b}(x)) : h_{w,b} \in L_d \}$$

is called the class of halfspace hypotheses.

Example in \mathbb{R}^2 :



$$w = (w_1, w_2)^T$$

$$\langle w, x \rangle + b = 0$$

$$w_1 \cdot x_1 + w_2 \cdot x_2 + b = 0$$

$$x_1 = 0$$

$$w_2 \cdot x_2 + b = 0$$

$$\Rightarrow x_2 = -\frac{b}{w_2}$$

Remark: Each affine function in \mathbb{R}^d can be re-written as a homogenous linear function in \mathbb{R}^{d+1} .

Ex.: Aff. function in \mathbb{R}^d : $\langle w, x \rangle + b$, $w \in \mathbb{R}^d$, $b \in \mathbb{R}$

Re-written as hom. lin. function in \mathbb{R}^{d+1} :

$$w' = (\underline{b}, w_1 \dots w_d)^T \text{ and}$$

$$x' = (\underline{1}, x_1 \dots x_d)^T$$

$$\Rightarrow \langle w, x \rangle + b = \langle w', x' \rangle$$

Question: What is the VC-dimension of halfspaces?

Theorem:

The VC-dimension, $VC(H)$, of the class of homogeneous halfspaces, i.e., $H = \{x \in \mathbb{R}^d : \langle w, x \rangle + b \geq 0\}$ is d !

Proof:

Part [1]: Lower-bound ($VC(H) \geq d$)

Let $C = \{x_1, \dots, x_d\}$ with $x_i \in \mathbb{R}^d$ and

$$x_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, x_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, x_d = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

If we set $w = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}$ for any given labeling (y_1, \dots, y_d) of our d data points

we get

$$\langle w, x_i \rangle = \sum_{j=1}^d w_j \cdot x_{ij} = y_i \cdot 1 = y_i$$

$\hookrightarrow x_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$ - i -th position

We conclude PART [1], i.e., $VC(H) \geq d$

PART [2] : Upper-bound ($V(H) < 0$)

We will establish $VC(H) \leq d+1$ via contradiction.

Assume $C = \{x_1, \dots, x_{d+1}\}$ is shattered by $H \Rightarrow$ All 2^{d+1} labelings can be realized! That means that we have weights w_k , $k=1, \dots, 2^{d+1}$ with which it is possible to realize the 2^{d+1} labelings.

Let's write:

$$\underbrace{\begin{pmatrix} w_1^T x_1 & w_2^T x_1 & \dots & w_{2^{d+1}}^T x_1 \\ w_1^T x_2 & w_2^T x_2 & & \\ \vdots & \vdots & & \\ w_1^T x_{d+1} & w_2^T x_{d+1} & \dots & w_{2^{d+1}}^T x_{d+1} \end{pmatrix}}_{2^{d+1}} = X \cdot w$$

X here is

$$X = \begin{pmatrix} -x_1 - \\ -x_2 - \\ \vdots \\ -x_{d+n} - \end{pmatrix} \Bigg| \underbrace{\quad}_{d}$$

W here is :

$$\left(\begin{array}{cccc|c} | & | & & | & \\ w_1 & w_2 & \dots & w_d & \vdots \\ | & | & & | & \end{array} \right) \xrightarrow{\text{d}}$$

$$(d+1) \times d \xrightarrow{X.W} d \times 2^{d+1}$$

Let $M = X \cdot w$. We see that $\text{sign}(M)$ gives All label combinations.

Claim: The rows of M are linearly independent. Why?

(Reminder: A sequence v_1, \dots, v_k in a vector space is lin. independent if

$$a_1 \cdot v_1 + a_2 \cdot v_2 + \dots + a_k \cdot v_k = 0$$

is only satisfied for $a_i = 0, i=1, \dots, k$)

! We see that the claim is true as there is no vector a with $a \neq 0$ such that $a^T M = 0$.

$$\underbrace{(-a^T -)}_{1 \times (d+1)} \cdot M = \underbrace{u}_{(d+1) \times (2^{d+1})} \in \mathbb{R}^{1 \times (2^{d+1})}$$

The k -th element in u is of the form $\boxed{a^T X w_k}$.

As we assume shattering, there is a k where

$$\text{sign}(a) = \text{sign}(X w_k)$$

As a result $a^T X w_k$ is a sum of positive numbers
 \Rightarrow this cannot be 0!

Lin. independence of the $d+1$ rows implies $\text{rank}(M) = d+1$.

We also know that

$$M = X \cdot W$$

and

$$\text{rank}(M) \leq \min(\text{rank}(X), \text{rank}(W)) \leq d$$

(a) only has d -columns

\Rightarrow CONTRADICTION! Our assumption that H shatters $d+1$ points does not hold!

This concludes PART 1 \square , i.e., $VC(H) < d+1$

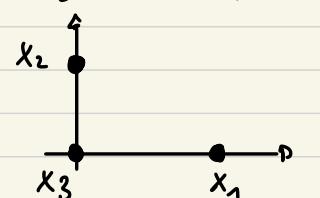
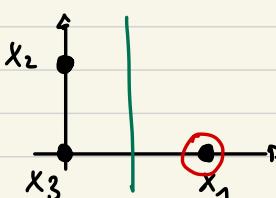
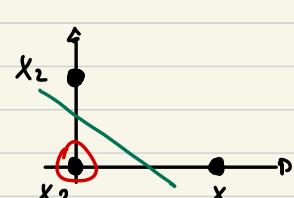
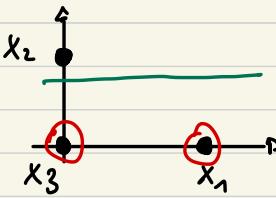
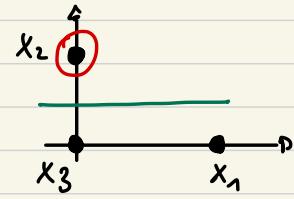
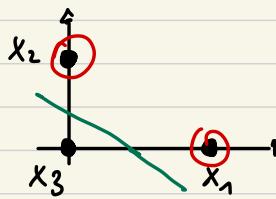
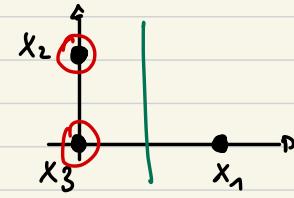
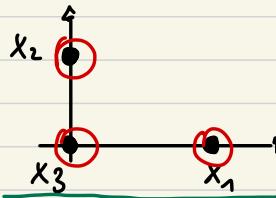
Overall, combining $VC(H) \geq d$ and $VC(H) < d+1$, we get

$$VC(H) = d$$



$$\langle \omega, x \rangle + b = \langle \omega', x' \rangle \quad x', \omega' \in \mathbb{R}^3$$

Example of halfspaces in \mathbb{R}^2 (VC-dim. lower bound)



x_1	x_2	x_3
1	1	1
1	1	-1
1	-1	1
1	-1	-1
-1	1	1
-1	1	-1
-1	-1	1
-1	-1	-1

$\rightarrow VC\text{-dim} \geq 3$

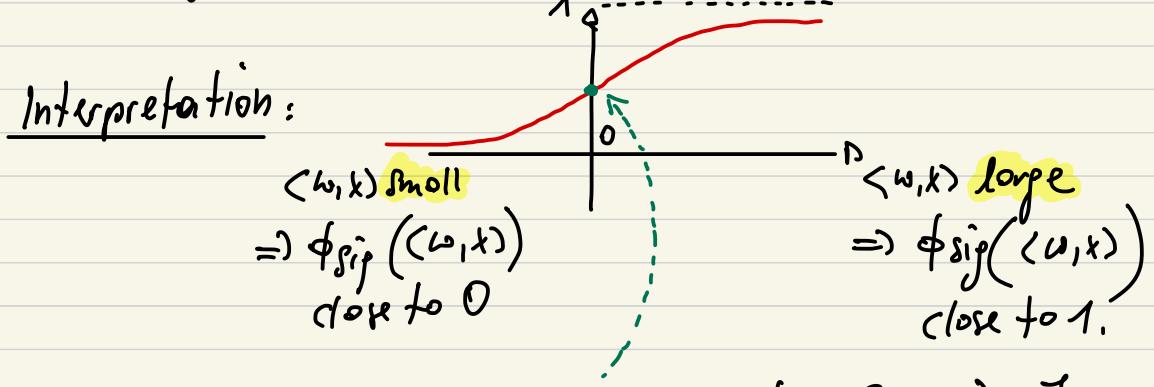
Linear predictors - Logistic Regression

Idea: we learn a class of functions from $\mathbb{R}^d \rightarrow [0, 1]$, i.e., $x \mapsto h(x) \in [0, 1]$ and interpret the output as the probability that x has label 1.

Formally, we compose L_d with ϕ_{sig} , defined as

$$\phi_{\text{sig}}(z) = \frac{1}{1 + e^{-z}} \quad (\text{Sigmoid function})$$

$$\Rightarrow H_{\text{sig}} = \underbrace{\phi_{\text{sig}} \circ L_d}_{\text{Hyp. class}} = \left\{ x \mapsto \phi_{\text{sig}}(\langle w, x \rangle) : w \in \mathbb{R}^d \right\}$$



Also, when $\langle w, x \rangle = 0 \Rightarrow \phi_{\text{sig}}(\langle w, x \rangle) = \frac{1}{2}$

Note that a halfspace hypothesis would always output 0 or 1 no matter if $\langle w, x \rangle$ is large/small or close to zero!
(since we have $\text{sign}(\langle w, x \rangle)$)

Since $h(x) \in [0, 1]$, we need to specify the "quality" of $h(x)$ given a training set $(x_1, y_1), \dots, (x_m, y_m)$ with, say $x_i \in \mathbb{R}^d$ and $y_i \in \{-1, +1\}$ \Rightarrow This leads us to the definition of a suitable loss function.

Intuition: if $y = +1$, we want $h(x)$ to be large.
 if $y = -1$, we want $h(x)$ to be small.
 (or, equivalently, $1 - h(x)$ to be large).

we have $1 - h(x) = 1 - \frac{1}{1 + e^{(\omega, x)}}$

$$= \frac{1 + e^{-(\omega, x)} - 1}{1 + e^{-(\omega, x)}}$$

$$= \frac{e^{-(\omega, x)}}{1 + e^{-(\omega, x)}}$$

$$= \frac{1}{e^{(\omega, x)} \cdot (1 + e^{-(\omega, x)})}$$

$$= \frac{1}{1 + e^{(\omega, x)}}$$

$\boxed{(x)}$
 if we write $\boxed{\frac{1}{1 + e^{-y(\omega, x)}}}$ and remember that in case of $y = +1$
 $h(x)$ should be large and if $y = -1$, then $1 - h(x)$ should be large

we realize that (x) plays a central role!

Any loss function that would make sense should thus monotonically increase with

$$\frac{1}{1 + e^{y \cdot \langle w, x \rangle}}$$

$$\boxed{\frac{1}{1 + e^{-y \cdot \langle w, x \rangle}}}$$

In logistic regression we take $\log\left(\frac{1}{1 + e^{-y \cdot \langle w, x \rangle}}\right)$ as loss, noting that $\log(\cdot)$ is a monotonic function.

\Rightarrow Given $(x_1, y_1), \dots, (x_m, y_m)$, ERM boils down to

$$(x) \quad \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \quad \frac{1}{m} \sum_{i=1}^m \log\left(\frac{1}{1 + e^{-y_i \cdot \langle w, x_i \rangle}}\right)$$

We will later show that this loss is convex in w , which will help solving (x) efficiently.

Next, we look at a probabilistic interpretation of this!

Logistic Regression - An alternative perspective

The idea is to consider the problem from a statistical perspective. Consider the conditional distribution of a response random variable Y , given an input random variable X , i.e.,

$$\underbrace{P[Y|X]}_{\text{tells us how precise our predictions are!}}$$

Now pick two classes, 0 or 1 \rightarrow then Y becomes an indicator variable

$$P[Y=1] = \mathbb{E}[Y]$$

$$\text{Also, } P[Y=1|X=x] = \mathbb{E}[Y|X=x]$$

Let's assume now that $P[Y=1|X=x] = p(x; \theta)$, for some function p that is parameterized by θ . In this case, the conditional likelihood function (under i.i.d. assumption) is

$$\prod_{i=1}^m P[Y=y_i | X=x_i] = \prod_{i=1}^m p(x_i; \theta)^{y_i} (1-p(x_i; \theta))^{1-y_i}$$

as y_1, \dots, y_m is a sequence of Bernoulli trials (where each trial has its own success probability - classically, we would have

$$\prod_{i=1}^m p^{y_i} (1-p)^{1-y_i} \quad \text{Now: } \prod_{i=1}^m p_i^{y_i} (1-p_i)^{1-y_i}$$

In the case

$$\prod_{i=1}^m p_i^{y_i} \cdot (1-p_i)^{1-y_i}$$

we would get (by max. likelihood) that $\hat{p}_i = 1$ if $y_i = 1$ and $\hat{p}_i = 0$ if $y_i = 0 \Rightarrow$ That does not make a lot of sense!

However, the p_i 's are not just arbitrary; in fact, they are linked together via $p_i = p(x_i; \theta)$. So p_i must be the same whenever the x_i 's are the same and if $p(\cdot; \cdot)$ is a continuous function, similar values of x_i give similar values of p_i .

We will assume that $p(\cdot; \cdot)$ is known, up to the parameters.
⇒ The likelihood function is a function of θ and we can use max. likelihood!

Choices for p : [1] linear function of x ? Not a good choice, as we need p to be in $[0, 1]$

[2] let $\log p$ be a linear function of x ?
Problem: $\log p$ is unbounded in one direction



13) "Logistic transform": let

$\log\left(\frac{p}{1-p}\right)$ a linear function of x !

\Rightarrow This means: $\log \frac{p(x; \theta)}{1-p(x; \theta)} = \langle \theta, x \rangle$

We get:

$$p(x; \theta) = \frac{e^{\langle \theta, x \rangle}}{1 + e^{-\langle \theta, x \rangle}} = \frac{1}{1 + e^{-\langle \theta, x \rangle}}$$

Using this, we can write down the log-likelihood function

$$\begin{aligned} & \sum_{i=1}^m y_i \cdot \log p(x_i; \theta) + (1-y_i) \cdot \log (1-p(x_i; \theta)) \\ &= \sum_{i=1}^m \log(1-p(x_i; \theta)) + y_i \cdot \log p(x_i; \theta) - y_i \cdot \log(1-p(x_i; \theta)) \\ &= \sum_{i=1}^m \log(1-p(x_i; \theta)) + y_i \cdot \log \left[\frac{p(x_i; \theta)}{1-p(x_i; \theta)} \right] \quad \left| \begin{array}{l} \langle \theta, x \rangle = \log\left(\frac{p}{1-p}\right) \\ (\text{from before}) \end{array} \right. \\ (x) &= \sum_{i=1}^m \log(1-p(x_i; \theta)) + y_i \langle \theta, x_i \rangle \end{aligned}$$

We know that $p(x; \theta) = \frac{1}{1+e^{-\langle \theta, x \rangle}} \Rightarrow 1-p(x; \theta) = \frac{1}{1+e^{\langle \theta, x \rangle}}$

$$(x) = \sum_{i=1}^m -\log(1+e^{\langle \theta, x_i \rangle}) + \underbrace{y_i \langle \theta, x_i \rangle}_{\log(e^{y_i \langle \theta, x_i \rangle})} = \sum_{i=1}^m \log\left(\frac{e^{y_i \langle \theta, x_i \rangle}}{1+e^{\langle \theta, x_i \rangle}}\right)$$

Now, simplifying gives

$$\sum_{i=1}^m \log \left[\frac{e^{y_i \cdot \langle \theta, x_i \rangle}}{1 + e^{\langle \theta, x_i \rangle}} \right] = \sum_{i=1}^m \log \left[\frac{1}{e^{y_i \cdot \langle \theta, x_i \rangle} + e^{(1-y_i) \cdot \langle \theta, x_i \rangle}} \right]$$

This would be the objective for max. likelihood. Equivalently, we can minimize the negative log-likelihood!

$$\Rightarrow \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left(\underbrace{e^{-y_i \cdot \langle \theta, x_i \rangle}}_{(1)} + e^{(1-y_i) \cdot \langle \theta, x_i \rangle} \right)$$

$$\text{For } y_i = 0 : (1) = 1 + e^{-\langle \theta, x_i \rangle}$$

$$\text{For } y_i = 1 : (1) = e^{-\langle \theta, x_i \rangle} + 1 = 1 + e^{-\langle \theta, x_i \rangle}$$

We see that this gives the same loss function as before, where we had $\mathcal{Y} = \{-1, +1\}$. Here, we have $\mathcal{Y} = \{0, 1\}$. But, the probabilistic perspective gives the same loss!

We also see that logistic regression gives a lin. classifier! But, the decision is more "fine-grain", as opposed to using $\text{sign}(\langle \omega, x \rangle)$ as in case of halfspaces.

