

①

Learning Infinite Classes II

In this lecture we introduce the notion of VCdim (Vapnik - Chernovenkis dimension) and discuss the fundamental theorem of learning theory.

An important result will be proved on the growth rate of hyp classes when the VCdim is finite (Sauer's lemma).

□.

Recap from last time.

$$C = \{c_1, c_2, \dots, c_m\} \subset X$$

\mathcal{H} set of hypothesis $h: X \rightarrow Y = \{0, 1\}$

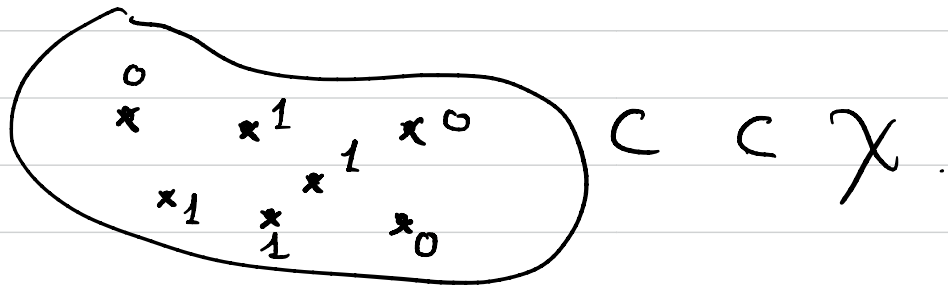
$\mathcal{H}_C =$ restriction of functions in \mathcal{H} to $C \rightarrow Y$

(2)

It is a good idea to view \mathcal{H}_C as the set of binary assignments

$$\{h(c_1) h(c_2) \dots h(c_m)\} \subset \{0,1\}^m$$

or labellings of elements of C .



Example: Threshold fct $\{h_a = \mathbb{1}_{x \leq a}\} = \mathcal{H}$



etc... \mathcal{H}_C here has 5 functions (5 possible labellings with threshold fct among all possible which are 16)

Definition of Shattering.

\mathcal{H} is said to shatter a set $C \subset X$ if

\mathcal{H}_C contains all possible labelings

or all possible fcts $C \rightarrow \{0, 1\}$ i.e. if $|\mathcal{H}_C| = 2^{|C|}$.

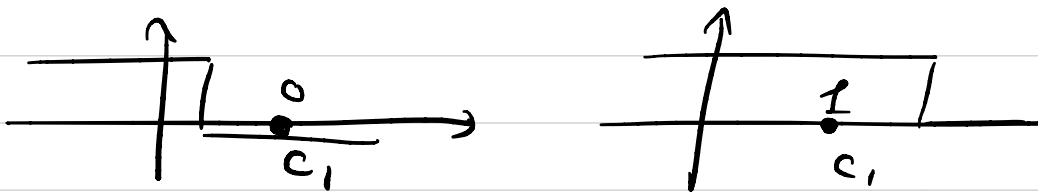
Examples.

- In the previous example above obviously

$\mathcal{H} = \text{Threshold fcts}$ does not shatter $C = \{c_1, c_2, c_3, c_4\}$.

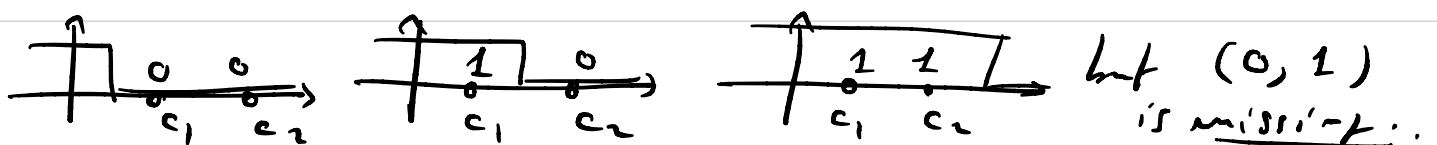
since we have $|\mathcal{H}_C| = 5$ and $2^{|C|} = 16$.

- $C = \{c_1\}$ Threshold fcts shatter C , indeed



so we get all possible labelings.

- $C = \{c_1, c_2\}$ Threshold fcts do not shatter C ,



(4)

Definition: VC dimension of an hypothesis class.

$VC \dim(\mathcal{H}) = d$ is the largest integer d such that we can find some set $C \subset \mathcal{X}$ with $|C| = d$ which is shattered by \mathcal{H} .

Remark: $VC \dim(\mathcal{H}) = d$ is the integer such that

- $\exists C$, $|C| = d$, shattered by \mathcal{H} so, s.t \mathcal{H}_C has all possible labelings.

AND:

- $\forall C$, $|C| \geq d+1$, are not shattered by \mathcal{H} so, s.t \mathcal{H}_C does not contain all possible labelings.

$$d = \max \{ |C| : \text{such that } |\mathcal{H}_C| = 2^{|C|} \}$$

Examples of computation of VC dim.

a) Threshold fcts $X = \mathbb{R}$. $\mathcal{H} = \{ \mathbb{1}_{x \leq a} = h_a \}$.

$C = \{c_1\}$ is shattered $\Rightarrow VCdim(\mathcal{H}) \geq 1$

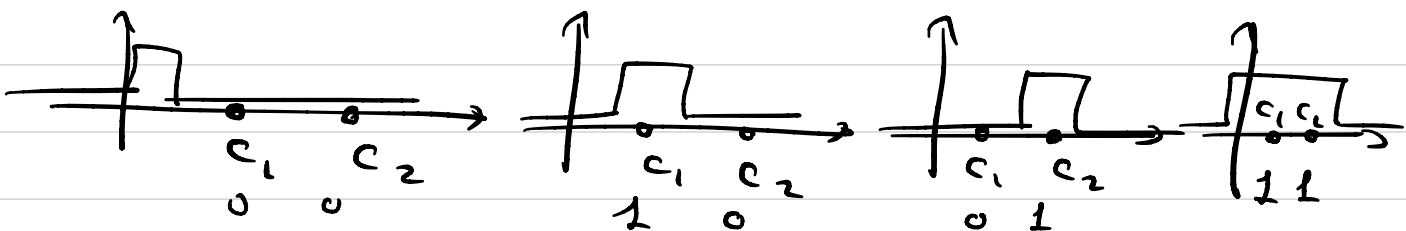
$C = \{c_1, c_2\}$ is not shattered $\Rightarrow VCdim(\mathcal{H}) < 2$

Thus $VCdim(\mathcal{H}) = 1$.

b) Interval fcts. $\mathcal{H} = \{ h_{a,b} = \mathbb{1}_{a \leq x \leq b} \}$

$C = \{c_1\}$ is shattered \rightarrow 

$C = \{c_1, c_2\}$ is also shattered



$\Rightarrow VCdim \mathcal{H} \geq 2$.

$C = \{c_1, c_2, c_3\}$ is not shattered. Indeed the labeling 101 is not obtained by rect fcts.

$\Rightarrow VCdim \mathcal{H} < 3$

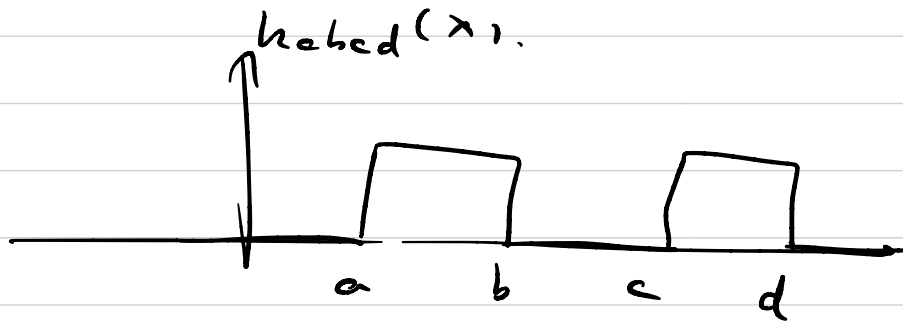
Thus $VCdim(\mathcal{H}) = 2$

6

c) Exercise;

$$\text{let } \mathcal{H}_{abcd} = \left\{ h_{abcd} = \mathbb{1}_{a \leq x \leq b} \cdot \mathbb{1}_{c \leq x \leq d} \right\}$$

for $a < b < c < d$.



Show that $\text{VCdim}(\mathcal{H}) = 4$.

d) Exercises.

$$\text{let } \mathcal{H} = \left\{ h_\theta : h_\theta(x) = \frac{1}{2} \lceil \sin \theta x \rceil, \theta \in \mathbb{R} \right\}$$

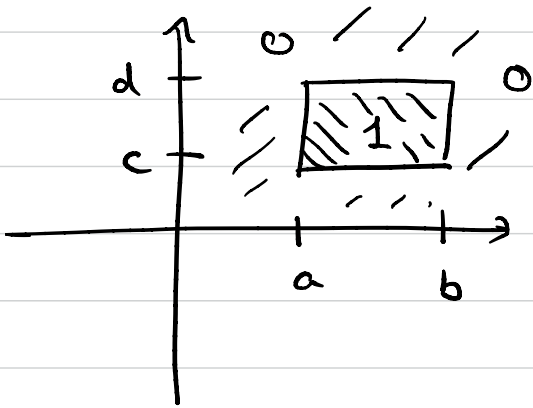
↑
upper integer part.

Note this class is parametrized by one parameter $\theta \in \mathbb{R}$.

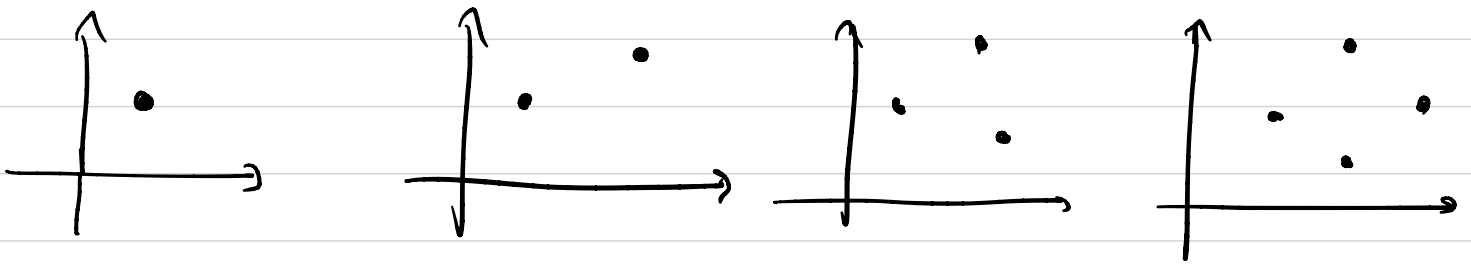
Show that $\text{VCdim}(\mathcal{H}) = \infty$. In other words given any m , one can construct C with $|C| = m$ s.t \mathcal{H}_C contains any labeling.

c) Axis aligned rectangles.

$$\chi = \mathbb{R}^2 \quad h_{[a,b] \times [c,d]}(x) = \begin{cases} 1 & x \in [a,b] \times [c,d] \\ 0 & \text{else} \end{cases}$$



fcn equal to 1 inside rect.



These sets C are all shattered \Rightarrow VCdim ≥ 4 .

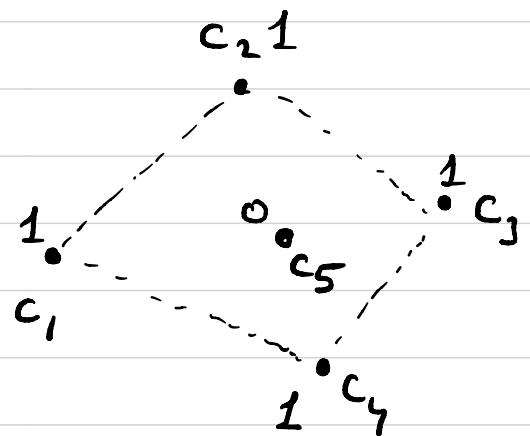
but All sets with $|C|=5$ cannot be shattered

$$c_1 = (c_1^x, c_1^y) \quad c_1^x \text{ min}$$

$$c_2 = (c_2^x, c_2^y) \quad c_2^y \text{ max}$$

$$c_3 = (c_3^x, c_3^y) \quad c_3^x \text{ max}$$

$$c_4 = (c_4^x, c_4^y) \quad c_4^y \text{ min}$$



VCdim < 5 .

$c_5 \rightarrow$ Necessarily inside convex hull

A Rectangle enclosing c_1, c_2, c_3, c_4 will also enclose c_5 \nearrow Hence.

The case of finite classes.

Consider a finite class $|\mathcal{H}| < +\infty$ and look at sets C such that $\underline{\underline{2^{|C|} > |\mathcal{H}|}}$.

Obviously these sets cannot be shattered since there are more potential labelings than hyp sets.

Thus sets C with $|C| > \log_2 |\mathcal{H}|$ cannot be shattered.

This means $VCdim(\mathcal{H}) \leq \log_2 |\mathcal{H}|$

(9)

No Free Lunch Theorem revisited.

Let $VCdim(\mathcal{H}) = +\infty$. This means for any integer $2m$, any set C of size $|C| = 2m$ is shattered and $|\mathcal{H}_C| = 2^{2m}$ i.e. C contains all possible labeling fcts. We can redo the proof of the No Free Lunch Theorem which was based on distr constructed out of all these labeling fcts. Conclusion was:

for any Alg $A(S)$ receiving $|S| \leq m \quad \exists \mathcal{D}_i$

over $\mathcal{X} \times \{0, 1\}$ and $f_i \in \mathcal{H}_C$ s.t

$$(a) \quad L_{\mathcal{D}_i}(f_i) = 0$$

$$(b) \quad \mathbb{P} \left(L_{\mathcal{D}_i}(A(S)) \geq \frac{1}{8} \right) \geq \frac{1}{7}$$

and hence \mathcal{H} is not PAC learnable.

Thm: $VCdim(\mathcal{H}) = +\infty \Rightarrow \mathcal{H}$ is not PAC learnable
 \mathcal{H} is PAC learnable $\Rightarrow VCdim(\mathcal{H}) < +\infty$.

Fundamental Theorem of PAC Learning.

Let $\mathcal{H} \ni h : \mathcal{X} \rightarrow \{0, 1\}$.

Let loss be a loss function s.t. $0 \leq \text{loss}(h, z) \leq 1$.

The following are all equivalent:

1. \mathcal{H} has the uniform conv property
2. \mathcal{H} is agnostic PAC learnable with the ERM rule
3. \mathcal{H} is agnostic PAC learnable.
4. \mathcal{H} is PAC learnable.
5. \mathcal{H} is PAC learnable with the ERM rule
6. \mathcal{H} is finite VC dimension.

Proof: $1 \Rightarrow 2$ (previous class)
 $2 \Rightarrow 3$, $3 \Rightarrow 4$, $2 \Rightarrow 5$ (trivial!)
 $4 \Rightarrow 6$ (No free lunch than just discussed).

$6 \Rightarrow 1$ Remains to be shown.

(11)

Now we discuss the proof of (6) \Rightarrow (1).

Main idea first:

Recap from last time $\tau_{\mathcal{H}_c}(m) = \max_{|C|=m} |\mathcal{H}_c|$

"growth rate"

and we proved

Thm: $\forall \delta$ we have for $\forall \epsilon > 0$

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left\{ \sup_{h \in \mathcal{H}} \left| \underbrace{L_{\mathcal{D}}(h)}_{\text{true Risk}} - \underbrace{L_S(h)}_{\text{empirical risk}} \right| \leq \frac{\epsilon + \sqrt{\log \tau_{\mathcal{H}_c}(m)}}{\delta \sqrt{2m}} \right\} \geq 1 - \delta$$

This then implies uniform convergence holds if we

have $\frac{\epsilon + \sqrt{\log \tau_{\mathcal{H}_c}(m)}}{\delta \sqrt{2m}} < \epsilon$ and thus \mathcal{H} is ϵ -quasi PAC learnable.

This can be achieved for $m \geq m_{\mathcal{H}_c}^{\text{UC}}(\epsilon, \delta) = C \frac{d + \log 1/\delta}{\epsilon^2}$

as long as $\tau_{\mathcal{H}_c}(m) = \text{poly}(m)$.

So what remains to be understood is the behavior of $\tau_{\mathcal{H}}(m)$ and notably when is it poly(m)?

- Note that if $m \leq d = \text{vc dim}(\mathcal{H})$ then

$\exists C$ with $|C| = m$ shattered by \mathcal{H} (by definition of VC)

$$\Rightarrow |\mathcal{H}_C| = 2^m \Rightarrow \boxed{\tau_{\mathcal{H}}(m) = 2^m \text{ for } m \leq d}$$

- The real question is what happens for $m \geq d+1$?

Lemma; Sauer-Shelah-Peres

Let $\text{vc dim}(\mathcal{H}) = d < +\infty$. Then for $m \geq d+1$

we have

$$\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d.$$

Proof of Lemma.

- The inequ $\sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d$ for $m \geq d+1$

is "calculus" and is left as exercise or (see Appendix of UDL).

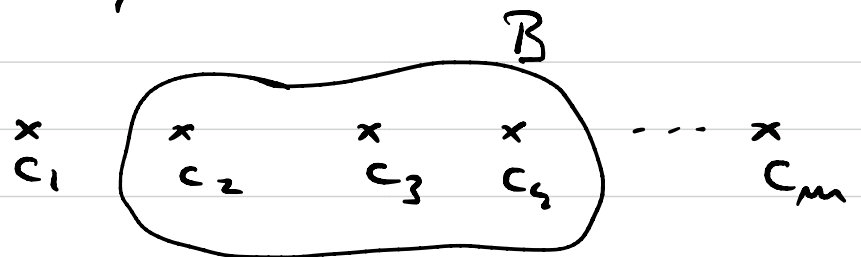
- We show here $\tau_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}$ for $m \geq d+1$.

First remark: it suffices to prove

$$|\mathcal{H}_C| \leq \left| \left\{ B \subset C \text{ such that } \mathcal{H} \text{ shatters } B \right\} \right|$$

Number of subsets of C that are shattered by \mathcal{H} .

Indeed:



B -subsets have size $|B| \leq d$ since they are shattered
 # of B -subsets of size i is $\binom{m}{i}$.

(Note for $i=0$ we take $B = \emptyset$ which is shattered by convention).

We will prove the inequ by induction.

$m=1$: $C = \{c\}$ possibilities for \mathcal{H}_C are

$$\mathcal{H}_C = \{h(c)=0\} \quad \mathcal{H}_C = \{h(c)=1\}$$

$$\text{and } \mathcal{H}_C = \{h_1(c)=0, h_2(c)=1\}.$$

- If $\mathcal{H}_C = \{h(c)=0\}$ or $\mathcal{H}_C = \{h(c)=1\}$

$$|\mathcal{H}_C| = 1 \text{ so } \tau_{\mathcal{H}}(m=1) = 1.$$

Subsets $B \subset C$ that are shattered is only \emptyset .
(since $C = \{c\}$ is not shattered here).

$$|\{B \subset C, B \text{ is shattered by } \mathcal{H}\}| = 1.$$

$$1 \leq 1 \quad \checkmark$$

- If $\mathcal{H}_C = \{h_1(c)=0, h_2(c)=1\}$ then

$$|\mathcal{H}_C| = 2 \text{ so } \tau_{\mathcal{H}}(m=2) = 2.$$

Subsets $B \subset C$ that are shattered are $\emptyset, C = \{c\}$

$$|\{B \subset C, B \text{ is shattered by } \mathcal{H}\}| = 2$$

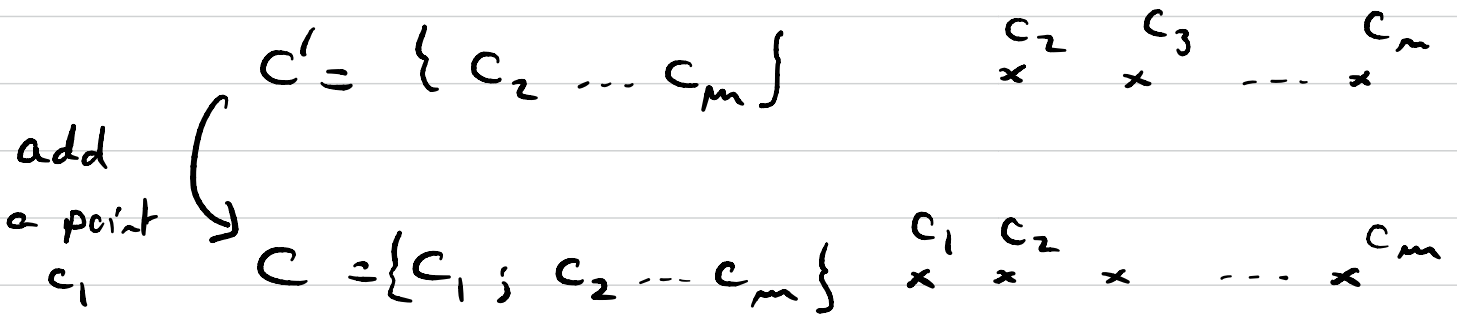
$$2 \leq 2 \quad \checkmark$$

Thus base case $m=1$ works.

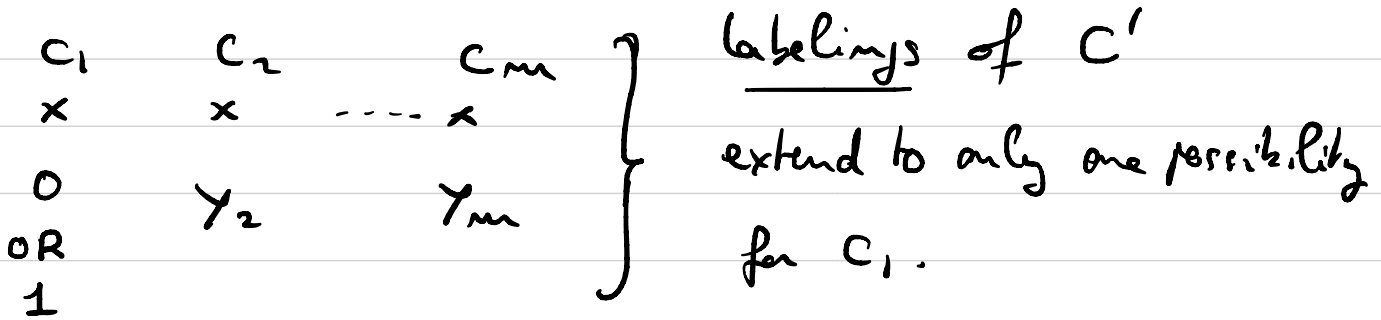
induction hypothesis: we assume the inequality

holds for all sets C' of size $1 \leq k \leq m-1$

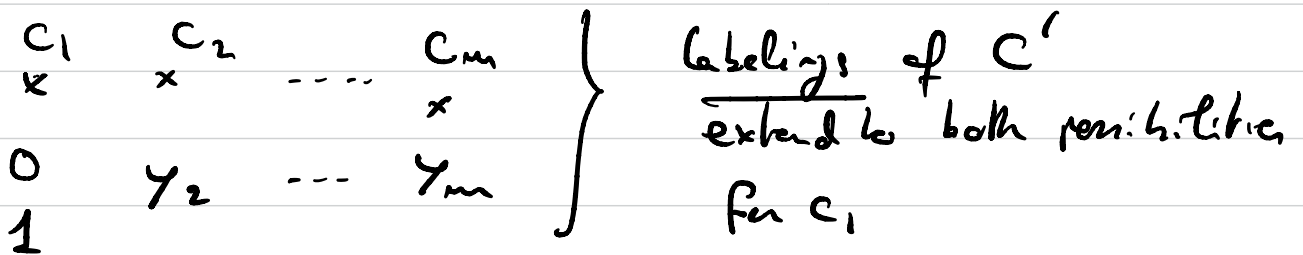
We must prove it then holds for C of size m .



Possibilities for H_C :



This set of labelings is called \mathcal{T}_0 .



This set of labelings is called \mathcal{T}_1 .

We have obviously

$$|\mathcal{H}_C| = |\mathcal{Y}_0| + |\mathcal{Y}_1|.$$

Inequality on $|\mathcal{Y}_0|$:

$$|\mathcal{Y}_0| = |\mathcal{H}_{C_1}|$$

$$\leq |\{B' \subset C' \text{ s.t. } B' \text{ is shattered by } \mathcal{H}\}|$$

↖ by induction hypothesis.

$$= |\{B \subset C \text{ s.t. } c_1 \notin B \text{ and } B \text{ is shattered by } \mathcal{H}\}|$$

↖ trivial

in summary

$$|\mathcal{Y}_0| \leq |\{B \subset C \text{ s.t. } c_1 \notin B \text{ and } B \text{ is shattered by } \mathcal{H}\}|$$

Inequality on Υ_1 :

Let $\tilde{\mathcal{H}} \subset \mathcal{H}$ where \mathcal{H} contains pairs of functions $\bar{h}, \bar{\bar{h}}$ which are equal on C' and opposite on C_1 .

$\tilde{\mathcal{H}} = \{ \bar{h} \in \mathcal{H} \text{ such that } \exists \bar{\bar{h}} \in \mathcal{H} \text{ with}$

$$(\underbrace{\bar{h}(c_1)}_{\text{m}}, \underbrace{\bar{h}(c_2)}_{\text{m}}, \dots, \underbrace{\bar{h}(c_m)}_{\text{m}}) = (\underbrace{1 - \bar{\bar{h}}(c_1)}_{\text{m}}, \underbrace{\bar{\bar{h}}(c_2)}_{\text{m}}, \dots, \underbrace{\bar{\bar{h}}(c_m)}_{\text{m}}) \}$$

agree

$$\bar{h}(c_1) \neq \bar{\bar{h}}(c_1)$$

$$|\Upsilon_1| = |\tilde{\mathcal{H}}_{c_1}|$$

$$\leq \left| \left\{ B' \subset C' \text{ such that } B' \text{ is shattered by } \tilde{\mathcal{H}} \right\} \right|$$

induction

$$= \left| \left\{ B' \subset C' \text{ s.t. } B' \cup \{c_1\} \text{ is shattered by } \tilde{\mathcal{H}} \right\} \right|$$

Since \mathcal{H} contains pairs of hyp which differ on C_1 .

$$= \left| \left\{ B \subset C \text{ such that } c_1 \in B \text{ and } B \text{ is trivial} \right. \right. \\ \left. \left. \text{shattered by } \tilde{\mathcal{H}} \right\} \right|$$

$$\leq \left| \left\{ B \subset C \text{ s.t. } c_1 \in B \text{ and } B \text{ is} \right. \right. \\ \left. \left. \text{shattered by } \mathcal{H} \right\} \right|$$

↑
because $\tilde{\mathcal{H}} \subset \mathcal{H}$ so more sets are shattered by \mathcal{H} since we have more fcts in \mathcal{H} .

in summary we have obtained

$$|\mathcal{Y}_0| \leq \left| \left\{ B \subset C \text{ s.t. } c_1 \notin B \text{ \& } B \text{ shattered by } \mathcal{H} \right\} \right|$$

$$|\mathcal{Y}_1| \leq \left| \left\{ B \subset C \text{ s.t. } c_1 \in B \text{ \& } B \text{ shattered by } \mathcal{H} \right\} \right|$$

$$\Rightarrow |\mathcal{H}_C| = |\mathcal{Y}_0| + |\mathcal{Y}_1|$$

$$\leq \left| \left\{ B \subset C \text{ s.t. } B \text{ shattered by } \mathcal{H} \right\} \right|$$

CQFD.