

PAC - Learning  
↓  
(Probably Approximately Correct)

Ref:  
FOMC book  
by Mohri et al.  
Chapters 2 & 3

The PAC Learning model:

$\mathcal{X}$ : input space (set of all feature vectors)

$\mathcal{Y}$ : set of labels e.g.,  $\{0, 1\}$

A concept  $c: \mathcal{X} \rightarrow \mathcal{Y}$

Concept class  $\mathcal{C}$ : collection of concepts

Suppose the inputs/examples are picked in an i.i.d. fashion using some distribution  $D$ .

The learning problem:

Given  $\mathcal{H}$ : hypothesis set (not necessarily  $= \mathcal{C}$ )

& data  $S = \{x_1, \dots, x_m\} \in \text{i.i.d. using } D$ ,

(labels  $= \{c(x_1), \dots, c(x_m)\}$ ),

the goal is to minimize the generalization error,

i.e.,  $R(h) = \mathbb{P}_{x \sim D} (h(x) \neq c(x)) = \mathbb{E} (1(h(x) \neq c(x)))$

given hypothesis  $h \in \mathcal{H}$

Empirical error:

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(h(x_i) \neq c(x_i))$$

By linearity of expectation,

$$\begin{aligned} E[\hat{R}_S(h)] &= \frac{1}{m} \sum_{i=1}^m E(\mathbb{1}(h(x_i) \neq c(x_i))) \\ &= \frac{1}{m} \sum_{i=1}^m E_{x \sim D}(\mathbb{1}(h(x) \neq c(x))) \\ &= E(\mathbb{1}(h(x) \neq c(x))) \\ &= R(h) \end{aligned}$$

fixed  
(not random)

### PAC-learning!

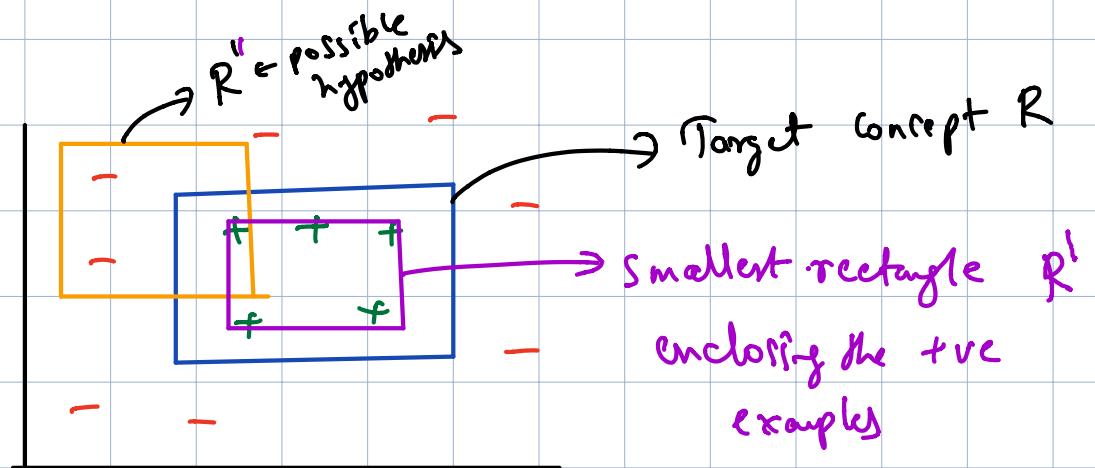
A concept class  $\mathcal{C}$  is PAC-learnable if there exists an algorithm  $A$  and a polynomial function  $\text{poly}(\cdot, \cdot, \cdot, \cdot)$  such that

$\forall \epsilon > 0, \delta > 0$ , for all distributions  $D$  on  $X$ ,  
and a target concept  $c \in \mathcal{C}$ , the  
following holds for any  $m \geq \text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, \text{size}(c), d)$ ,

$$P_{S \sim D^m} (R(h_S(A)) \leq \epsilon) \geq 1 - \delta.$$

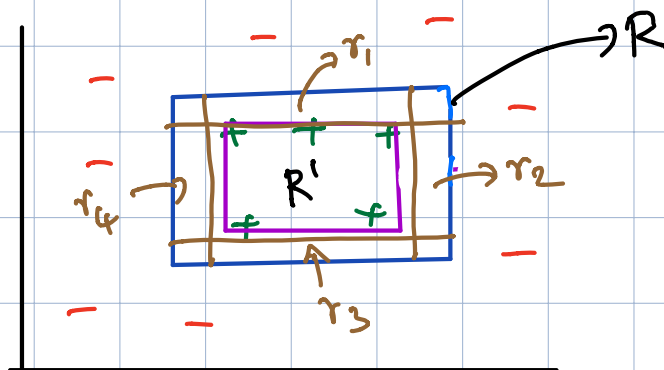
- Note
- (i) PAC model is distribution-free
  - (ii) Training & test samples come from the same distribution
  - (iii) Learnability  $\rightarrow$   $\mathcal{H}$  of the concept class

Example: (Medium-build or axis-aligned rectangles)



Fix  $\epsilon > 0$ . Let  $P(R)$  be the probability mass of  $\mathcal{R}$

Suppose that  $P(R) > \epsilon$ .



$$R = [l, r] \times [b, t], \quad r_4 = [l, s_4] \times [b, t],$$

where  $s_4 = \min \{s \mid P([l, s] \times [b, t]) \geq \epsilon/4\}$

If  $R'$  has one side in each  $\sigma_i$ , then it's error (which is the prob. of region not covered by  $R'$ ) is  $< \epsilon$ .

Now, If  $P(R') > \epsilon$ , then  $R'$  misses at least one of the regions.

Given data  $S = \{x_1, \dots, x_m\}$ , let  $R_S$  be the smallest rectangle enclosing positive examples.

Then,

$$P(R(R_S) > \epsilon) \leq P\left(\bigcup_{i=1}^4 \{R_S \cap \sigma_i = \emptyset\}\right)$$

$$\stackrel{\text{union bound}}{\leq} \sum_{i=1}^4 P(\{R_S \cap \sigma_i = \emptyset\})$$

$$\stackrel{(\text{since } p(\sigma_i) \geq \epsilon/4)}{\leq} 4 \left(1 - \frac{\epsilon}{4}\right)^m$$

$$\stackrel{(\text{used } 1-x \leq e^{-x})}{\leq} 4 \exp\left(-\frac{m\epsilon}{4}\right)$$

To ensure  $P(R(R_S) > \epsilon) \leq \delta$ ,

it is enough if we have

$$4 \exp\left(-\frac{m\epsilon}{4}\right) \leq \delta$$

$\uparrow$

$$m \geq \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right)$$



Thus,  $\forall \epsilon > 0, 0 < \delta < 1$ , if  $m \geq \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right)$ ,

then  $P(R(h_S) > \epsilon) < \delta$ .

$$\text{or } P(R(h_S) \leq \epsilon) \geq 1 - \delta$$

## Guarantees for finite hypothesis sets - consistent case

Theorem: Let  $\mathcal{H}$  be a finite set of functions  $h: \mathcal{X} \rightarrow \mathcal{Y}$ . Let  $A$  be an algorithm that for any target concept  $c \in \mathcal{H}$ , and training data  $S$ , returns a "consistent" hypothesis  $h_S$ , i.e., empirical error  $\rightarrow \hat{R}_S(h_S) = \frac{1}{m} \sum_{i=1}^m 1(h_S(x_i) \neq c(x_i)) = 0$

Then, for any  $\epsilon, \delta > 0$ ,

$$P(R(h_S) \leq \epsilon) \geq 1 - \delta \text{ holds if}$$

$$m \geq \frac{1}{\epsilon} \left( \log |\mathcal{H}| + \log \frac{1}{\delta} \right).$$

Proof: fix  $\epsilon > 0$ . "Uniform Convergence".

Note:  $h_S$  is random as it depends on  $S = \{x_1, \dots, x_m\}$

We bound the probability that some  $h \in \mathcal{H}$  is consistent and has  $R(h) > \epsilon$

$$\text{Let } \mathcal{H}_\epsilon = \{h \in \mathcal{H} \mid R(h) > \epsilon\}$$

Given Dataset  $S$ ,

$$P(\hat{R}_S(h)=0) \leq (1-\epsilon)^m$$

$\uparrow$   
"fixed  $h$ "

$$\begin{aligned} &P(\exists h \in H_\epsilon : \hat{R}_S(h)=0) \\ &= P(\hat{R}_S(h_1)=0 \text{ (or) } \hat{R}_S(h_2)=0 \dots \text{(or) } \hat{R}_S(h_{|H_\epsilon|})=0) \\ &\leq \sum_{h \in H_\epsilon} P(\hat{R}_S(h)=0) \\ &\leq \sum_{h \in H_\epsilon} (1-\epsilon)^m \\ &\leq |H| (1-\epsilon)^m \\ &\leq |H| \exp(-m\epsilon) \quad \leftarrow \text{equate this to } \delta \text{ \& find an expression for } m \text{ in terms of } |H|, \epsilon, \delta. \\ &\leq \delta \quad \text{if} \\ &\quad m \geq \frac{1}{\epsilon} \left( \log |H| + \log\left(\frac{1}{\delta}\right) \right) \end{aligned}$$

So, we have

$$P(R(h_S) > \epsilon) \leq \delta$$

$$\text{or } P(R(h_S) \leq \epsilon) > 1-\delta$$



Equivalent statement:

$$\text{Set } \delta = |H| \exp(-m\epsilon)$$

& find an expression for  $\epsilon$ , i.e.,

$$\epsilon = \frac{1}{m} \left( \log |H| + \log \frac{1}{\delta} \right)$$

$$P(R(h_s) \leq \frac{1}{m} \left( \log |H| + \log \frac{1}{\delta} \right)) > 1 - \delta$$

Example: Conjunction of Boolean literals (at most "2" of them)

Boolean literal:  $x_i$  or  $\bar{x}_i$

Conjunction:  $x_1 \wedge \bar{x}_2 \wedge x_4$

0	1	1	0	1	1	+
0	1	1	1	1	1	+
0	0	1	1	0	1	-
0	1	1	1	1	1	+
1	0	0	1	1	0	-
0	1	0	0	1	1	+
0	1	?	?	1	1	

$$\Rightarrow \bar{x}_1 \wedge x_2 \wedge x_5 \wedge x_6$$

Given a dataset, can we find a consistent hypothesis?

Note: negative examples are not informative.  
so, focus on the positive examples.

Algorithm: Start with all literals, say  
 $x_1 \wedge \bar{x}_1, x_2 \wedge \bar{x}_2, \dots, x_d \wedge \bar{x}_d$

& rule out literals that are incompatible with positive examples.

$$|H| = 3^d$$

Using the bound, we get the PAC guarantee

$$\text{for } m \geq \frac{1}{\epsilon} \left( d \log 3 + \log \frac{1}{\delta} \right)$$

E.g. for  $\delta = 0.02$ ,  $\epsilon = 0.1$ ,  $d = 10$ ,

the bound is  $m \geq 149$ , i.e.,

with at least 149 samples, we can obtain a hypothesis that is 90% accurate with probability at least 98%.

---

Guarantees for a finite hypothesis set, where consistency is not available.

Probability bounds: Hoeffding's inequality

Hoeffding's lemma:

A. r.v.  $X \in [a, b]$ ,  $E X = \mu$ ,  $b \neq a$

Then,

$$E(\exp^{\lambda(X-\mu)}) \leq \exp\left(\frac{\lambda^2(b-a)^2}{8}\right)$$

Proof  $\leftarrow$  skipped.

Hoeffding's inequality:

Let  $X_1, \dots, X_m$  be iid samples, with  $X_i \in [a_i, b_i]$ . Then, letting  $S_m = \sum_{i=1}^m X_i$ ,

$$P(S_m - E(S_m) > t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

Special case:  $X_i \in [a, b]$

$$P(S_m - E(S_m) > t) \leq \exp\left(-\frac{2t^2}{m(b-a)^2}\right)$$

$$\bar{X}_m = \frac{1}{m} S_m, \quad E X_i = \mu$$

$$P(\bar{X}_m - \mu > \epsilon)$$

$$= P\left(\frac{S_m}{m} - \mu > \epsilon\right) = P(S_m - m\mu > m\epsilon)$$

$$\leq \exp\left(-\frac{2 m \epsilon^2}{(b-a)^2}\right)$$

with Chebyshev's inequality,  
the RHS is  $\frac{(b-a)^2}{m\epsilon^2}$

## Proof of Hoeffding's inequality:

For any r.v.  $X$  &  $\lambda > 0$ , we have <sup>to be specified later.</sup>

$$P(X \geq t) = P(\exp(\lambda X) \geq \exp(\lambda t))$$

Markov  
ineq.

$$\leq \frac{E(\exp(\lambda X))}{\exp(\lambda t)}$$

$$P(S_m - ES_m \geq t)$$

$$\leq e^{-\lambda t} E(\exp(\lambda(S_m - ES_m)))$$

Samples are  
independent  $\rightarrow$

$$\leq e^{-\lambda t} \prod_{i=1}^m E(\exp(\lambda(X_i - EX_i)))$$

Hoeffding's  
lemma  $\rightarrow$

$$\leq e^{-\lambda t} \prod_{i=1}^m \exp\left(\frac{\lambda^2 (b_i - a_i)^2}{8}\right)$$

$$= e^{-\lambda t} \exp\left(\lambda^2 \sum_{i=1}^m \frac{(b_i - a_i)^2}{8}\right)$$

Set  $\lambda = \frac{4t}{\sum_{i=1}^m (b_i - a_i)^2}$

$$\leq \exp\left(-\frac{2t^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

Using a parallel argument, one can obtain

$$P(S_m - ES_m \leq -\epsilon) \leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

Combining, we have

$$P(|S_m - ES_m| > \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^m (b_i - a_i)^2}\right)$$

---

Back to PAC-learning with finite hypothesis sets -  
the "inconsistent" case

Using Hoeffding's inequality, for any  
"fixed" hypothesis  $h: X \rightarrow \{0, 1\}$ , we have

$$P(\hat{R}_S(h) - R(h) \geq \epsilon) \leq \exp(-2m\epsilon^2)$$

over the  
distribution of  $S$   
that has " $m$ " iid  
samples

$$P(\hat{R}_S(h) - R(h) \leq -\epsilon) \leq \exp(-2m\epsilon^2)$$

(or)

$$P(|\hat{R}_S(h) - R(h)| \geq \epsilon) \leq 2 \exp(-2m\epsilon^2)$$

Example: (Tossing a coin)

A coin with bias  $p$  (=prob(heads)) is tossed " $m$ " times.  
iid

hypothesis  $h =$  predict tails.

$$R(h) = p, \quad \hat{R}_S(h) = \hat{p} \rightarrow \text{sample avg of heads}$$

$$P(|\hat{p} - p| \geq \epsilon) \leq 2 \exp(-2m\epsilon^2)$$

$\downarrow$   
equate to  $\delta$

$$P\left(|\hat{p} - p| \geq \sqrt{\frac{\log 2/\delta}{2m}}\right) \leq \delta$$

$$P\left(|\hat{p} - p| \leq \sqrt{\frac{\log 2/\delta}{2m}}\right) \geq 1 - \delta$$

E.g.  $\delta = 0.02$ ,  $m = 500$ . Then,

With probability at least 98%,

$$|\hat{p} - p| \leq \sqrt{\frac{\log(100)}{1000}} \approx 0.067$$

Claim: Let  $\mathcal{H}$  be a finite hypothesis set.

Then, for any  $\delta > 0$ , w.p.  $(1 - \delta)$ , we have

$$\forall h \in \mathcal{H}, \quad R(h) \leq \hat{R}_S(h) + \sqrt{\frac{\log |\mathcal{H}| + \log 2/\delta}{2m}}$$

$S$  is a dataset with " $m$ " iid samples

Remark: In comparison to the bound in the consistent case, we have a bound on that scales with  $\sqrt{\frac{\log |\mathcal{H}|}{m}}$  in the inconsistent case.



In the consistent case, the bound scaled as  $\frac{\log |H|}{n}$ .

Proof: Let  $h_1, \dots, h_{|H|}$  be the elements of  $H$ .

Then,

$$P(\exists h \in H \text{ s.t. } |\hat{R}_S(h) - R(h)| > \epsilon)$$

$$= P(|\hat{R}_S(h_1) - R(h_1)| > \epsilon \text{ (or) } \\ |\hat{R}_S(h_2) - R(h_2)| > \epsilon \text{ (or) } \\ \vdots$$

$$|\hat{R}_S(h_{|H|}) - R(h_{|H|})| > \epsilon)$$

$$\leq \sum_{h \in H} P(|\hat{R}_S(h) - R(h)| > \epsilon)$$

$$\leq 2 |H| \exp(-2n\epsilon^2)$$

To complete the proof, set

$$\delta = 2 |H| \exp(-2n\epsilon^2)$$



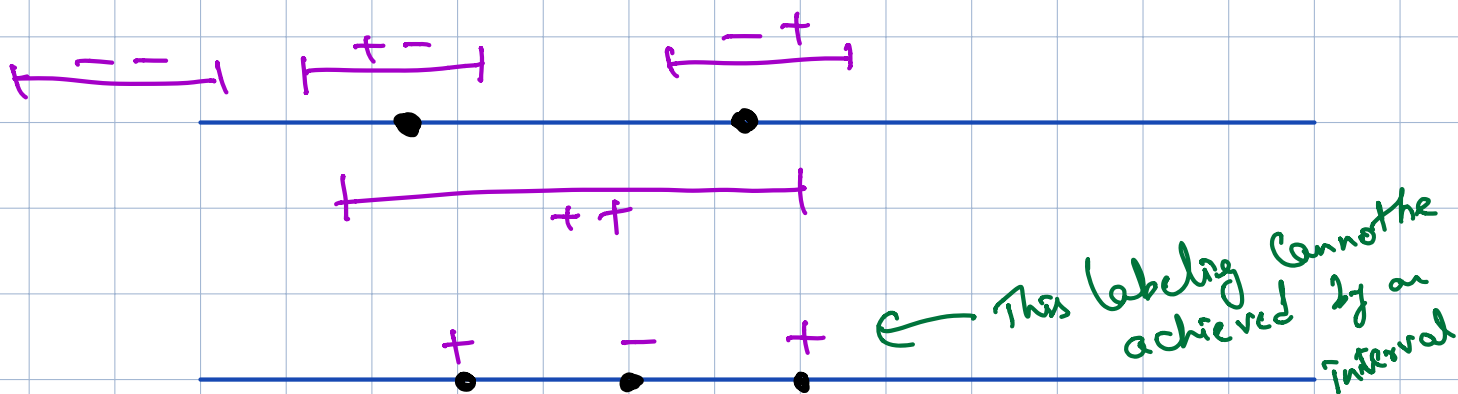
# The Case of infinite hypothesis sets:

## Growth function & VC-dimension

Growth function  $\Pi_H : \mathbb{N} \rightarrow \mathbb{N}$  for a hypothesis set  $H$  is defined as

$$\forall m \in \mathbb{N}, \Pi_H(m) = \max_{\{x_1, \dots, x_m\} \subseteq X} |\{h(x_1), \dots, h(x_m) \mid h \in H\}|$$

Example: Intervals on a real line



Claim: "Growth function generalization bound"

Let  $H$  be a family of functions taking values in  $\{-1, +1\}$ . Then, for any  $\delta > 0$ , w.p.  $(1-\delta)$

$$R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2 \log \Pi_H(m)}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

## VC - dimension:

Shattering? A set  $S$  of  $m$  points is said to be shattered by a hypothesis set  $H$  if

$$\Pi_H(m) = 2^m$$

Def: The VC-dimension of a hypothesis set  $H$  is the size of the largest set that can be shattered by  $H$ .

$$VCdim(H) = \max \{ m \mid \Pi_H(m) = 2^m \}$$

### Remark:

① If  $VCdim(H) = d$ , then there exists a set of size  $d$  that can be shattered by  $H$ .

This does not imply that all sets of size " $d$  or less than  $d$ " are shattered by  $H$ .

② To give a lower bound for  $\text{VCdim}(\mathcal{H})$ , it is enough to show a set of size  $d$  that can be shattered by  $\mathcal{H}$  i.e.,

$$\text{VCdim}(\mathcal{H}) \geq d$$

For the upper bound, one has to show that no set  $S$  of cardinality " $d+1$ " can be shattered by  $\mathcal{H}$ , i.e.,

$$\text{VCdim}(\mathcal{H}) < d+1$$

Example 1: "Intervals on the real line"

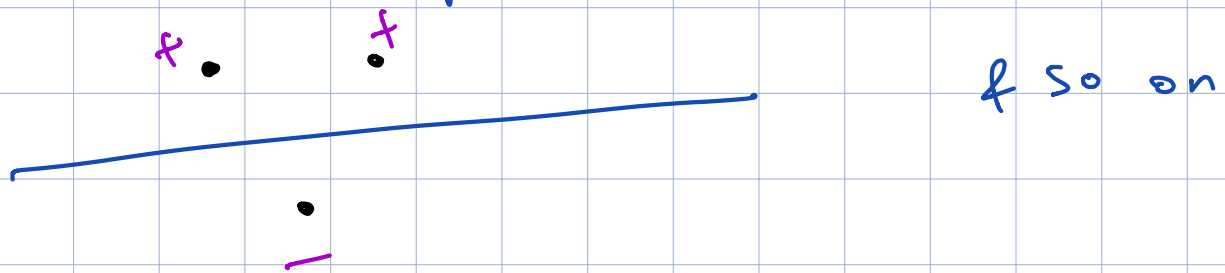
$$\text{VCdim}(\text{"set of intervals"}) = 2.$$

Example 2: "Hyperplanes".

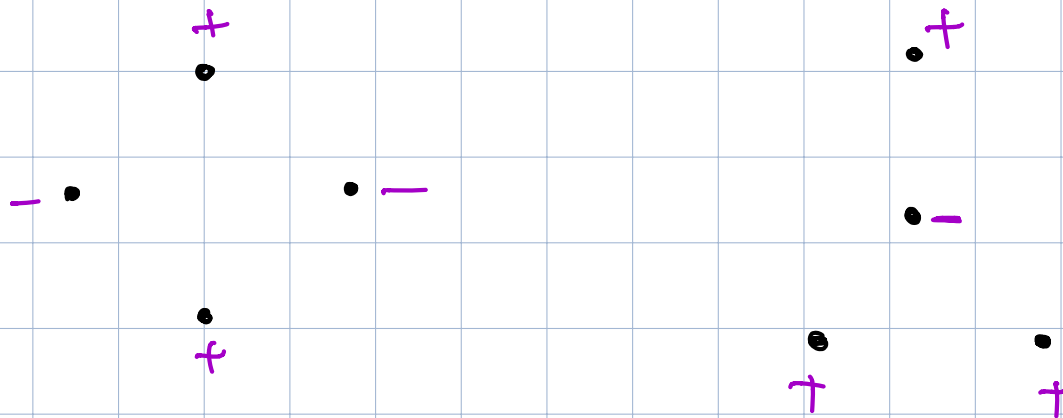
Consider the set of hyperplanes in  $\mathbb{R}^2$

$$\text{VCdim}(\text{hyperplanes in } \mathbb{R}^2) \geq 3$$

Set  $S =$  "3 non-collinear points"



To show  $VC(\text{dim}(\text{hyperplanes})) < 4$



So,  $VC(\text{dim}(\text{hyperplanes})) = 3$ .

This result can be generalized to  $\mathbb{R}^d$ , i.e.,

$$VC(\text{dim}(\text{hyperplanes in } \mathbb{R}^d)) = d + 1$$

Lower bound: Need a set  $S$  of " $d+1$ " points that can be shattered.

$$S = \{x_0, x_1, \dots, x_d\} \quad x_i \in \mathbb{R}^d$$

$$x_0 = \text{Origin} = (0, \dots, 0)$$

$$x_1 = e_1 = (1, 0, \dots, 0)$$

$$\vdots$$

$$x_d = e_d = (0, \dots, 0, 1)$$

$y_0, \dots, y_d \in \{+1, -1\}$  be the labeling

$$w = (y_1, \dots, y_d)$$

hyperplane for classification:

$$w^T x + \frac{y_0}{2}, \text{ where } x \in \mathbb{R}^d$$

$$\text{Want: } \text{sgn} \left( w^T x_i + \frac{y_0}{2} \right) = y_i$$

The above holds because

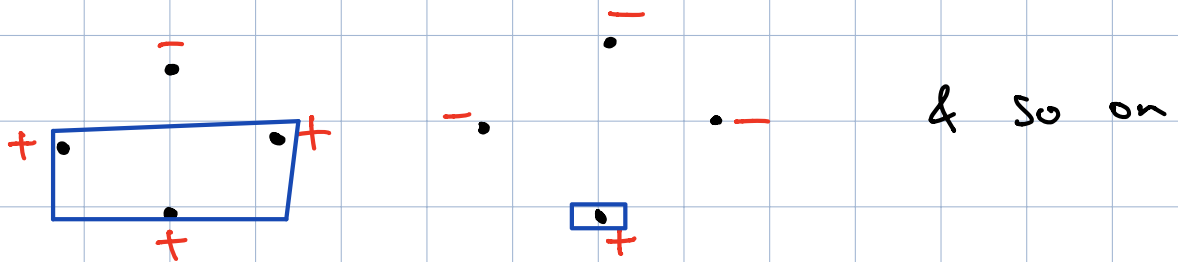
$$\begin{aligned} & \text{sgn} \left( w^T x_i + \frac{y_0}{2} \right) \\ &= \text{sgn} \left( y_i + \frac{y_0}{2} \right) = y_i \quad \text{for } i=1, \dots, d \end{aligned}$$

One can show an upper bound of  $d+2$  for  $VCdim(\text{hyperplanes in } \mathbb{R}^d)$  using "Radon's theorem" (check the textbook).

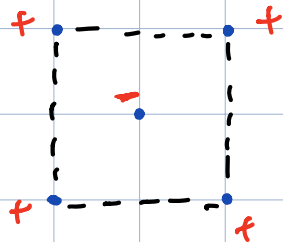
$$VCdim(\text{hyperplanes in } \mathbb{R}^d) = d+1$$

Example 3: Axis-aligned Rectangles (AAR)

$$VCdim(AAR) \geq 4$$



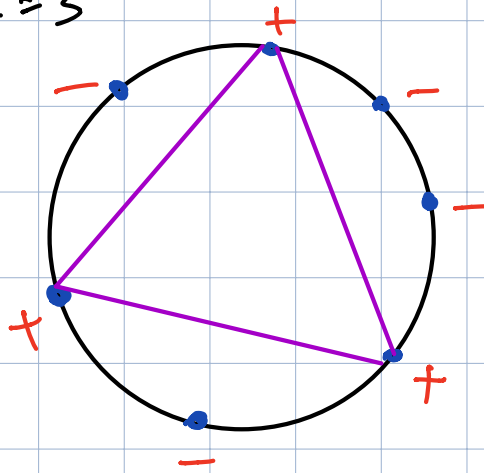
$$VCdim(AAR) < 5$$



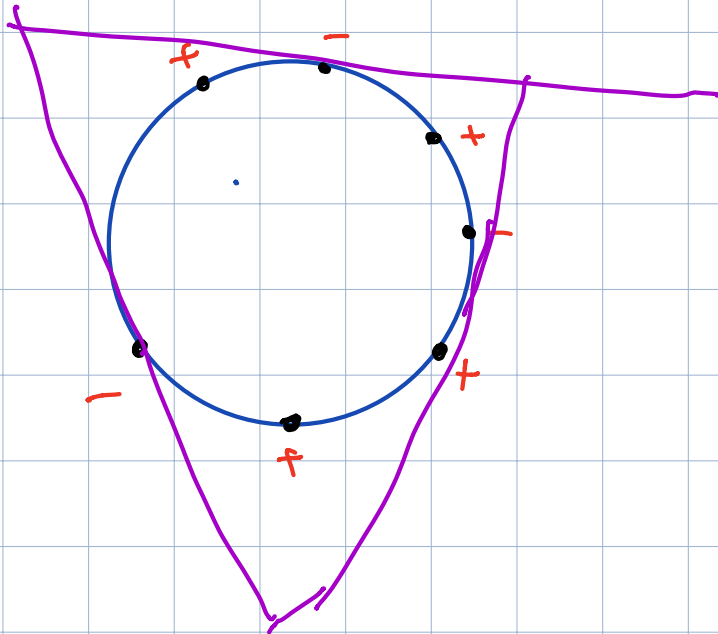
Example: "Convex  $d$ -gons in the plane"

Claim:  $VCdim(\text{convex } d\text{-gons}) \geq 2d+1$

Ex.  $d=3$



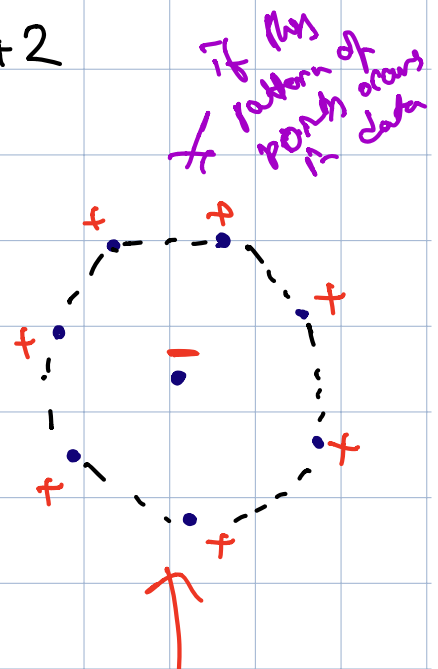
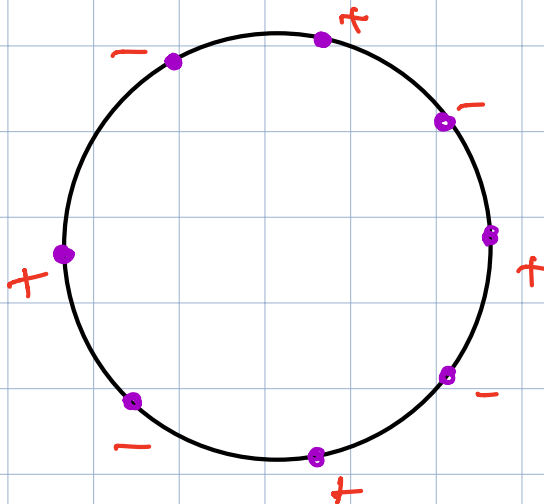
$\hookrightarrow \# +ves < \# -ves$



$\hookrightarrow \# +ves > \# -ves$

Upper bound

$$V(dim(d\text{-gons})) < 2d+2$$





↑  
Case where no point  
is in the convex  
hull of remaining  
points.

a point is in the  
interior of the rest  
of the points  
(= convex hull of  
the remaining points)

Hence,  $VCdim(d\text{-gons in the plane}) = 2d + 1$

Note:  $VCdim(\text{convex polygons}) = +\infty$ .

$$(or) \quad \prod_{\text{convex polygons}}(n) = 2^n \quad \forall n \geq 1$$

(Sauer's lemma): Let  $H$  be a hypothesis set  
with  $VCdim(H) = d$ . Then,  $\forall n \geq 1$ ,

$$\prod_H(n) \leq \sum_{i=0}^d \binom{n}{i}$$

"No proof".

Corollary: Let  $H$  be a hypothesis set with  $VCdim(H) = d$ .

Then,  $\forall n \geq d$ ,

$$\prod_H(n) \leq \left(\frac{en}{d}\right)^d$$

Pf:-

From Sauer's lemma,

$$\begin{aligned}\Pi_{\mathcal{H}}(m) &\leq \sum_{i=0}^d \binom{m}{i} \\ &\leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \\ &\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \\ &= \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i \\ &= \left(\frac{m}{d}\right)^d \left(1 + \frac{d}{m}\right)^m \\ &\leq \left(\frac{m}{d}\right)^d e^d\end{aligned}$$

---

Note:- Two cases  $\rightarrow V(\dim(\mathcal{H})) = d < \infty \Rightarrow \Pi_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$   
 $\rightarrow V(\dim(\mathcal{H})) = \infty \Rightarrow \Pi_{\mathcal{H}}(m) = 2^m$

---

Main claim  $\Leftrightarrow$  Generalization bound with  
VC-dimension

Let  $\mathcal{H}$  be a family of functions taking values  
in  $\{-1, +1\}$  &  $V(\dim(\mathcal{H})) = d$ .

For any  $\delta > 0$ , w.p.  $(1-\delta)$ ,  $\forall h \in \mathcal{H}$ ,

$$\mathcal{R}(h) \leq \hat{\mathcal{R}}_S(h) + \sqrt{\frac{2d \log \frac{em}{\delta}}{m}} + \sqrt{\frac{\log 1/\delta}{2m}}$$