


Consistency,

Uniform Convergence,

and Learning:

The VC Dimension

Finite H : Quick Review

- Any consistent \hat{h} has

$$\mathbb{E}(\hat{h}) \leq \mathbb{E} \quad \text{as long as } \checkmark$$
$$m > \frac{1}{\epsilon} \log \frac{12H}{\delta}$$

- Even allows $\|\hat{h}\|_2 > \sqrt{m}$

- $\frac{1}{\epsilon} \rightarrow \frac{1}{\epsilon^2}$ buys us

uniform convergence:

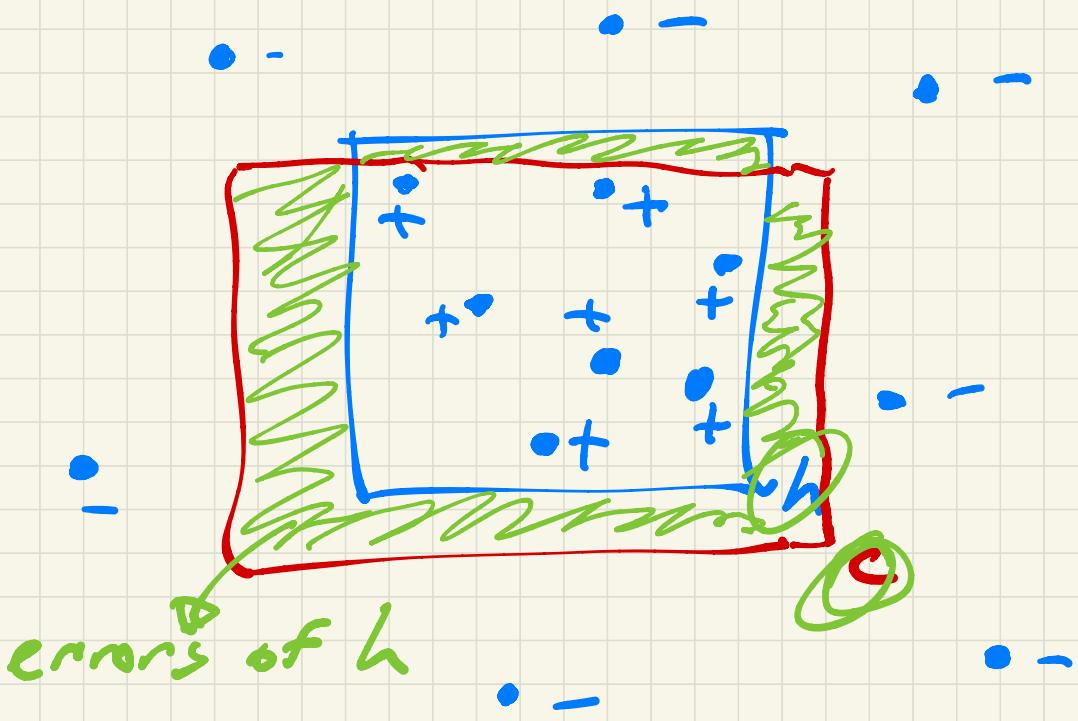
$$|\mathbb{E}_{\hat{h}}(h) - \mathbb{E}(h)| \leq \epsilon$$

$\forall h \in \mathcal{H}$

What happens when
 $|z_1|$ is infinite?

What quantity should
replace $\log|z_1|$?

What's wrong with this:



If $\epsilon(h) \geq \epsilon$, then

prob we missed it

$$\leq (1-\epsilon)^n \leq \delta$$

• Let $S = \{x_1, x_2, \dots, x_m\} \subseteq X$
be an ordered set
(note: no labels)

Key definition #1:

$$\Pi_{\mathcal{H}}(S) \triangleq \{ (h(x_1), h(x_2), \dots, h(x_m)) : h \in \mathcal{H} \}$$

\mathcal{H}
 $\subset \{0, 1\}^m$

Set of all labelings of
 S induced by \mathcal{H} .

Key definition #2:

Say that S is shattered by \mathcal{H} if $\Pi_{\mathcal{H}}(S) = \{0, 1\}^m$

Key definition #3:

The VC dimension of \mathcal{H} is the size of the largest set shattered by \mathcal{H} .

Denote by $VC(\mathcal{H})$.

Quantification

To show $VC(H) \geq d$,
we must find some
shattered S of size d .

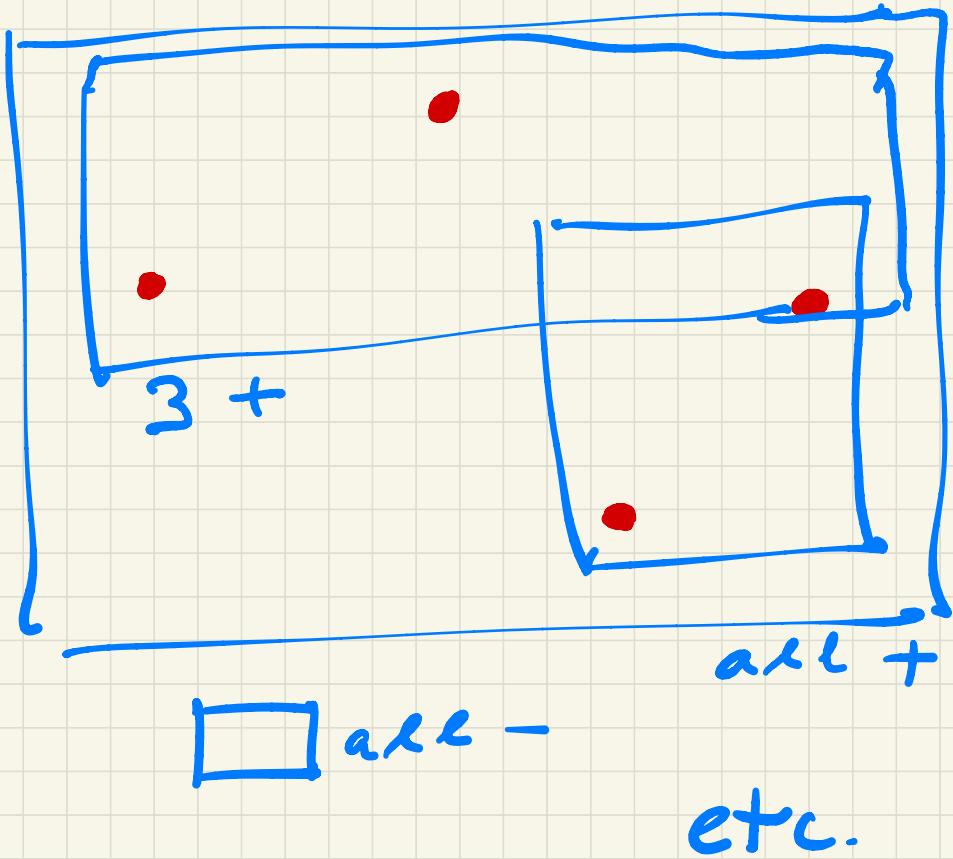
To show $VC(H) \leq d$,
we must show that
no set of size $d+1$
is shattered.

Roadmap

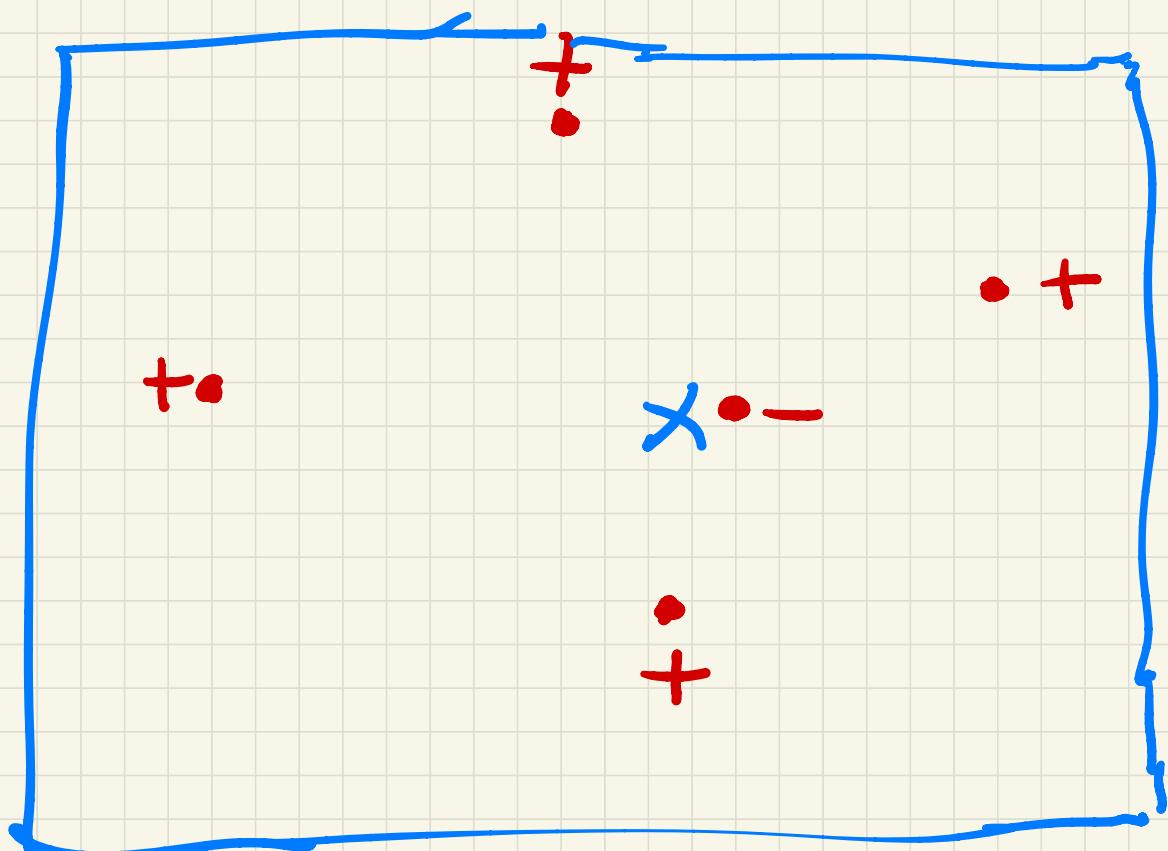
- We will see that $VC(H)$ replaces $\log |H|$ in proving consistency \Rightarrow PAC
- Warm-up with ex's of competing $VC(H)$
- Main result has both combinatorial and probabilistic components

Ex: Rectangles in \mathbb{R}^2

Can shatter these 4 points:



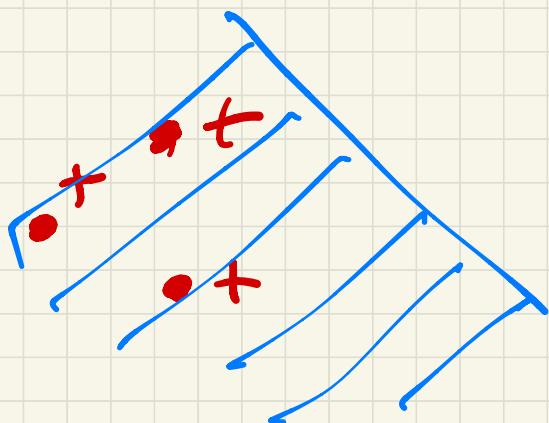
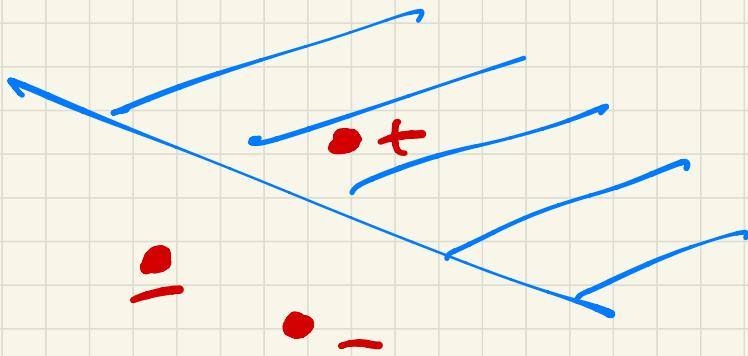
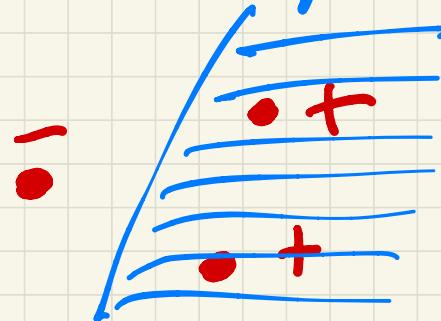
Can shatter no 5 points:



$$\therefore \text{VC dim} = 4$$

Ex: Halfspaces in \mathbb{R}^2

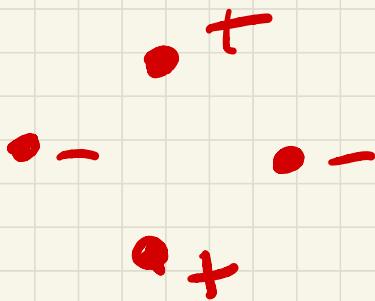
Can shatter 3 points -



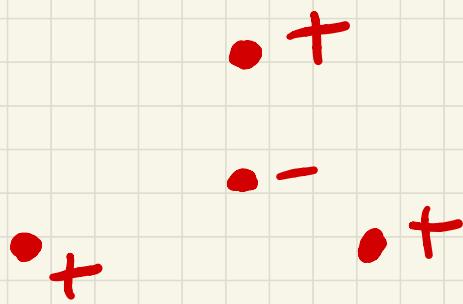
etc.

Can shatter no 4 points:

E.g.



and

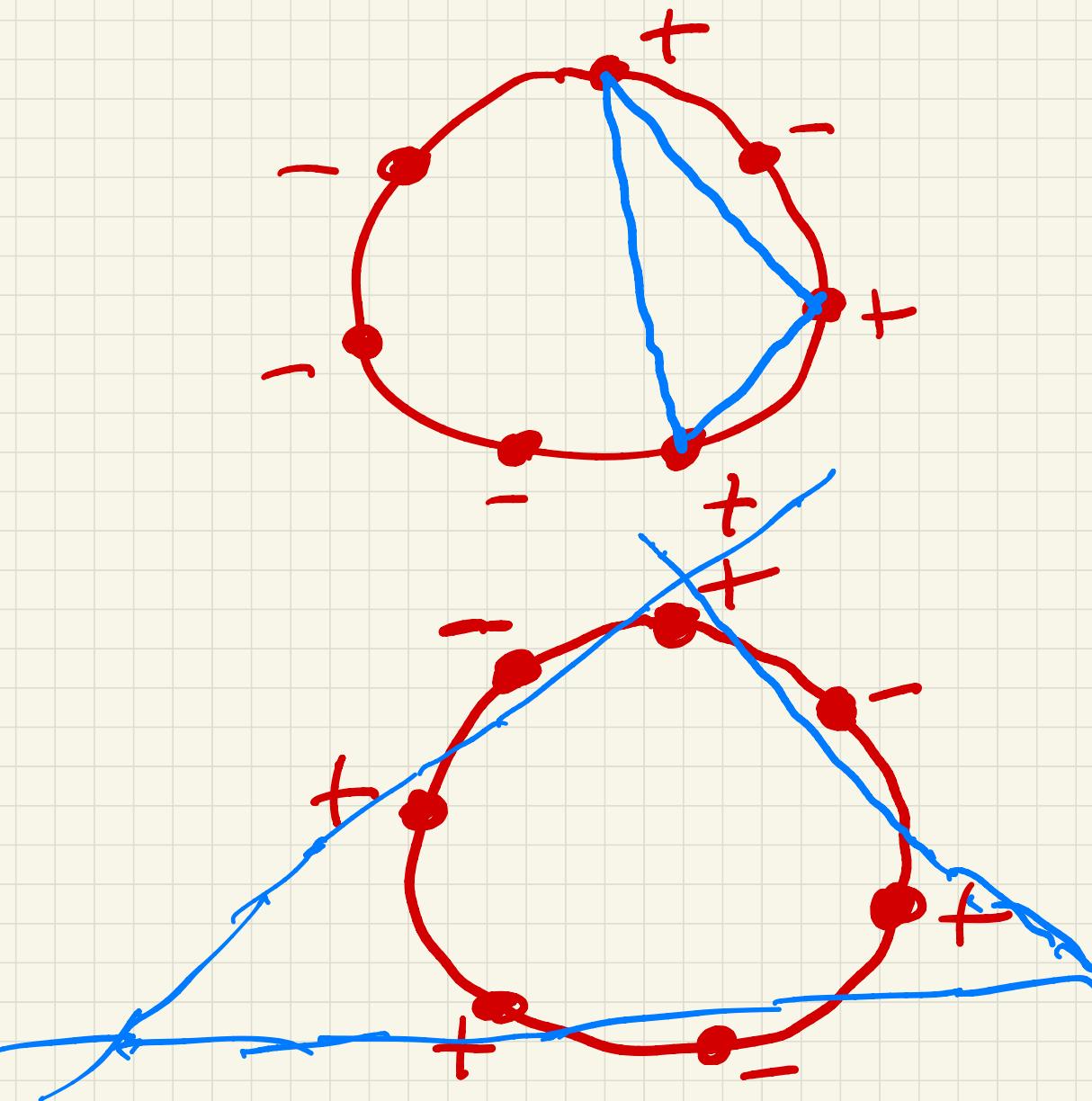


VC dim = 3

In \mathbb{R}^d : VC dim = $d+1$

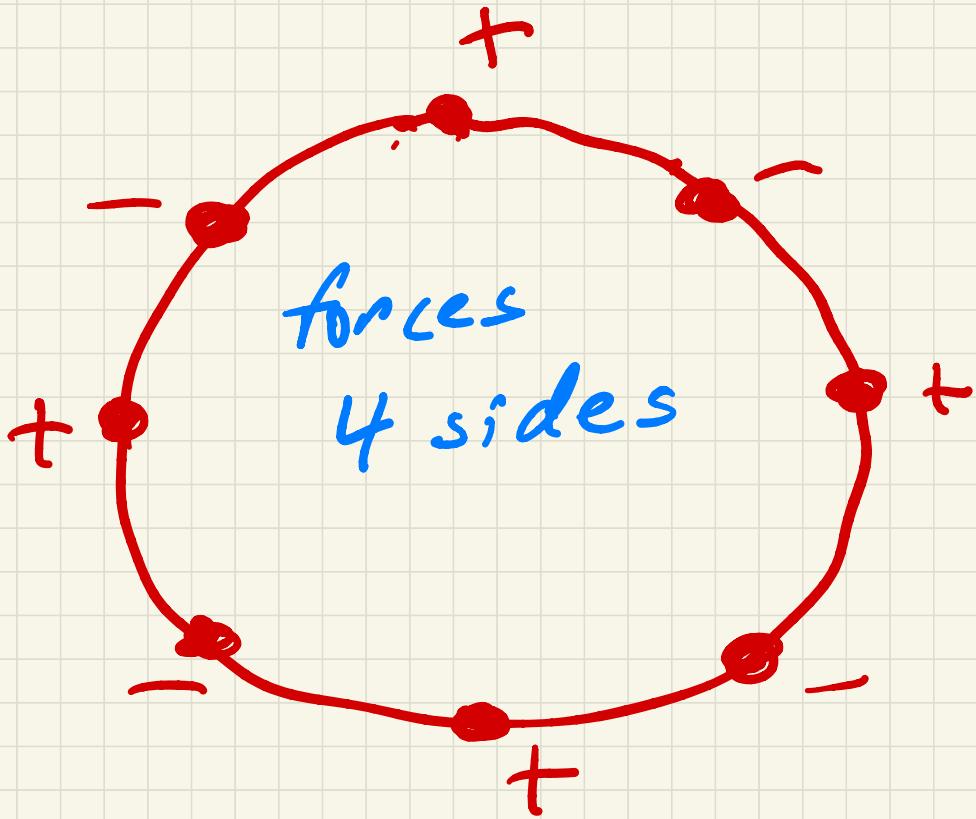
Ex: d-gons in \mathbb{R}^2

Can shatter $2d+1$ pts:



Can shatter no $2d+1$ pts:

Worst case:



$$VC \dim = 2d + 1$$

Ex: Conjunctions over $\{0, 1\}^n$

Can shatter these n pts:

0111 ...

-

10111 ...

+

11011 ...

-

11101 ...

-

111101 ...

+

$x_1 x_3 x_4 \dots$

all -
indices

Note: For any finite \mathcal{H} ,

$$\text{VC dim} \leq \log_2 |\mathcal{H}|$$

since shattering d pts
requires 2^d functions.

So VC dim conjunctions

$$\leq \log 3^n$$

$$\therefore \text{VC dim} = \Theta(n)$$

Recall

$\Pi_{\mathcal{H}}(S)$ = set of labelings
of S induced by \mathcal{H}

Let's define

$$\Pi_{\mathcal{H}}(m) = \max_{S: |S|=m} \{ |\Pi_{\mathcal{H}}(S)| \}$$

"growth function"

Then for $m \leq \text{VC}(\mathcal{H})$:

$$\Pi_{\mathcal{H}}(m) = 2^m \text{ (max)}$$

What about $m > \text{VC}(\mathcal{H})$?

Amazing fact:

Let $d = \text{VC}(\mathcal{H})$. Then

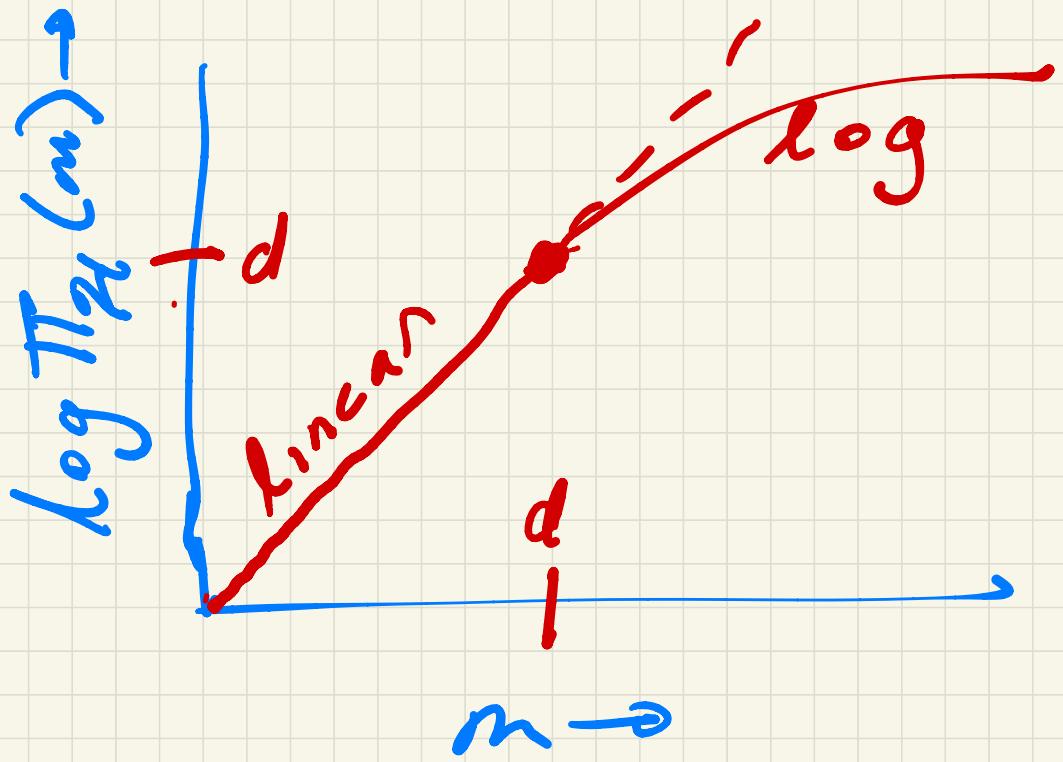
for $m > d$, $\text{Tr}_{\mathcal{H}}(m) \sim m^d$.

$m \leq d$: exponential growth

$m > d$: polynomial growth

"Sauer's Lemma"

Visually:



Overall, $T_{lg}(m) = O(m^d)$

$$(2^m \leq m^m)$$

Proof outline

1. Define a function $\phi_d(m)$, show that
 $T\Theta_1(m) \leq \phi_d(m)$
2. Show that
 $\phi_d(m) = O(m^d)$

Define:

$$\phi_d(m) = \phi_d(m-1) + \phi_{d-1}(m-1),$$

$$\phi_0(m) = \phi_d(0) = 1.$$

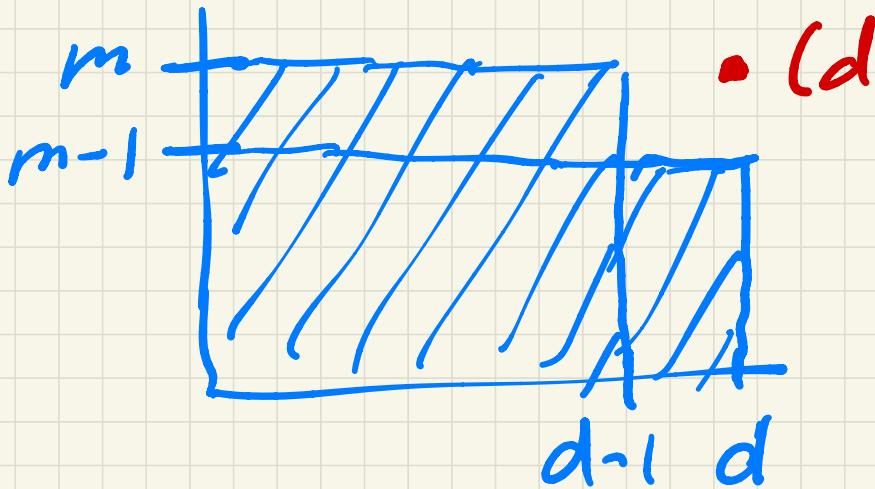
Lemma: If $VC(\mathcal{H}) = d$,
then fim

$$\text{fim}_{\mathcal{H}}(m) \leq \phi_d(m)$$

Let's prove by

double induction
on d & m .

Assume true for any
 $d' \leq d$ (one must)
 $m' \leq m$ (be strict)



$S:$

x_1	x_2	x_3	\dots	x_{m-1}	x_m
0	0	1	\dots	0	1
1	1	1	\dots	1	1
0	0	1	\dots	0	1
			\vdots		
1	1	0	\dots	0	0

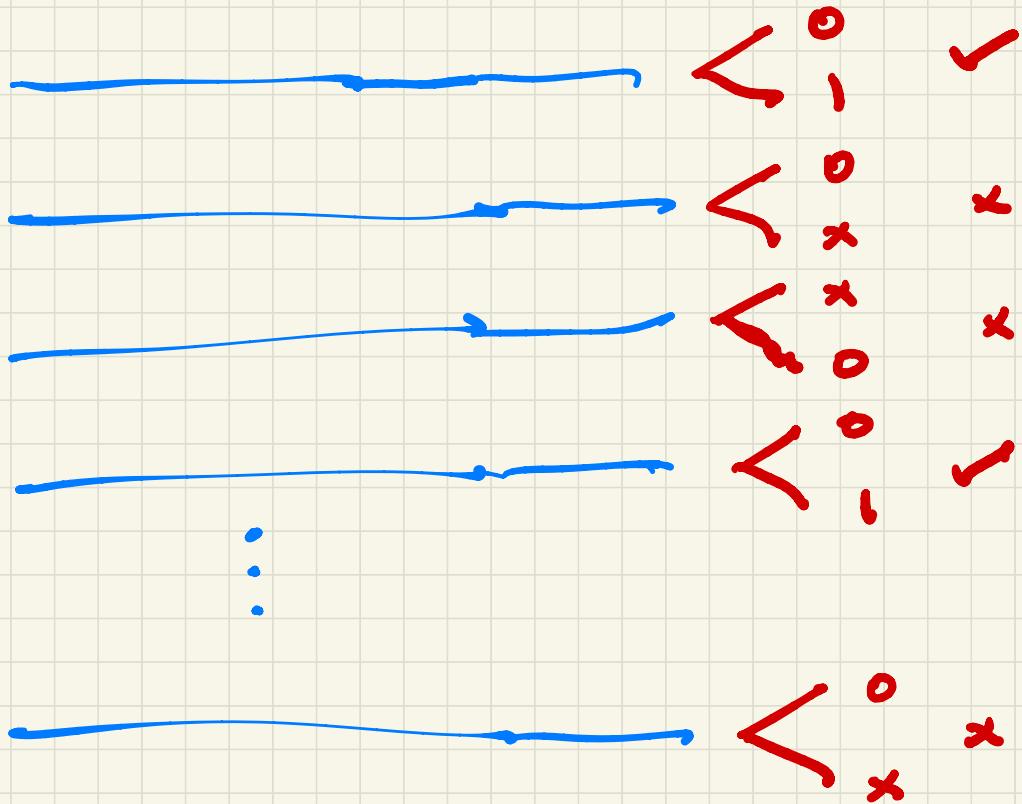
How many labelings here?

$$\leq \pi_{\mathcal{H}}(m-1) \leq \phi_d(m-1)$$

(induction)

labelings on
 x_1, \dots, x_{m-1} :

extension(s)
on x_m :



Need to add

count of ✓'s.

But labellings of

$$X' = x_1, x_2, \dots, x_{m-1}$$

with a ✓ is just a fn.
class (over X')

of $\text{VC dim} \leq d-1$.

Why?

$$\therefore \#\sqrt{s} \leq \Phi_{d-1}(m-1)$$

(induction again)

$$\# \Pi_H(s) \leq \Phi_d(m-1)$$

$$+ \Phi_{d-1}(m-1)$$

$$= \Phi_d(m)$$

as desired.

$$\underline{\text{Lemma}} \quad \Phi_d(m) = \sum_{i=0}^d \binom{m}{i}.$$

Pf $\Phi_d(m) =$

$$\Phi_{d-1}(m) + \Phi_{d-1}(m-1) =$$

$$\sum_{i=0}^d \binom{m-1}{i} + \sum_{i=0}^{d-1} \binom{m-1}{i}$$

(induction)

$$\xrightarrow{\text{hi}} \binom{m-1}{d-1}$$

$$= \sum_{i=0}^d \binom{m-1}{i-1} \xrightarrow{\text{lo}} \binom{m-1}{-1} \triangleq 0$$

$$= \sum_{i=0}^d \left[\binom{m-1}{i} + \binom{m-1}{i-1} \right] = \binom{m}{i}$$

OK, that's it for

combinatorial
analysis of $\pi_H(m)$.

Now for the
probabilistic
part.

- Let's fix all of these:
target $c \in \mathcal{H}, P, \epsilon, \delta$
- Define error regions
of \mathcal{H} wrt c :

$$\mathcal{H}_{\Delta c} \triangleq \left\{ \underbrace{h_{\Delta c}}_{\text{viewed as sets of } + \text{ex's}} : h \in \mathcal{H} \right\}$$

As functions:

$$(h_{\Delta c})(x) = \begin{cases} 1 & \text{if } h(x) \neq c(x) \\ 0 & \text{else} \end{cases}$$

$$= h(x) \oplus c(x)$$

Note: $VC(\mathcal{H}_{\Delta C}) = VC(\mathcal{H})$:

• If $h_1(x) \neq h_2(x)$,

$$h_1(x) \oplus c(x) \neq h_2(x) \oplus c(x)$$

• labelings $h_1(s) \neq h_2(s)$

$$\Rightarrow (h_1 \oplus c)(s) \neq (h_2 \oplus c)(s)$$

$$\therefore \forall s | \pi_{\mathcal{H}}(s) | = | \pi_{\mathcal{H}_{\Delta C}}(s) |$$

$$\Rightarrow VC(\mathcal{H}) = VC(\mathcal{H}_{\Delta C})$$

OK. Suppose some $h \in \mathcal{H}$
has $\varepsilon(h) > \varepsilon$. We want to
make sure S hits
 $h \in \mathcal{H}$ at least once.

And we want this to happen
for every h s.t. $\varepsilon(h) > \varepsilon$.

Call such S an ε -net
and define

$$A(s) = \begin{cases} 1 & \text{if } S \text{ not an } \varepsilon\text{-net} \\ 0 & \text{else} \end{cases}$$

Goal: bound $\Pr_S[A(s)=1]$

Problem: Hard to tell if $A(s)=1$ by looking at S .

Solution: "two-sample trick"

$B(s, s') = 1$ iff $\exists r \in \mathcal{H} \Delta C$
2m pts s.t. $P[r] > \epsilon$
and:

① $r \cap S = \emptyset \ (\Rightarrow A(s) = 1)$

② $|r \cap S'| \geq \frac{\epsilon m}{2}$

Otherwise $B(s, s') = 0$.

Claim: $\Pr[B] \geq \Pr[A] \cdot \frac{1}{2}$

$B(s, s') = 1 \Rightarrow A(s) = 1$, so

just need to lower bound

$$\Pr_{S, S'} [B(s, s') = 1 \mid A(s) = 1]$$

can fix missed r

s.t. $P[r] > \epsilon$

$$\geq \Pr_{S'} [1 \cap r s''] \geq \epsilon m/2$$

$\geq \frac{1}{2}$ (e.g. Chebyshev)

$$\text{So } \Pr[A] \leq \underline{2\Pr[B]}$$

upper bound

Draw S, S' in 2 steps:

- draw $2m$ pts $T = S \cup S'$
At most $\Phi_d(2m)$ labelings!
- split $T \rightarrow S, S'$ randomly
(random permutation)

Same as i.i.d. S, S' by
exchangeability

$B(S, S') = 1$ only if there are
s.t. r is hit $\ell \geq \epsilon m / 2$

times in T but all ℓ
hits fall in S'

Prob. of above happening
just wrt
random
permutation

$$\leq \frac{\binom{m}{\ell}}{\binom{2m}{\ell}}$$

$$\frac{\binom{m}{\ell}}{\binom{2m}{\ell}} = \frac{\frac{m! \ell!}{(m-\ell)!}}{\frac{(2m)! \ell!}{(2m-\ell)!}}$$

$$= \frac{m!}{(m-\ell)!} \frac{(2m-\ell)!}{(2m)!}$$

$$= \frac{m(m-1)(m-2)\dots(m-\ell+1)}{(2m)(2m-1)\dots(2m-\ell+1)}$$

$$= \frac{m}{2m} \cdot \frac{m-1}{2m-1} \cdot \frac{m-2}{2m-2} \cdots \frac{m-\ell+1}{2m-\ell+1}$$

$$\leq 2^{-\ell} \leq 2^{-\varepsilon m/2}$$

Wrapping up:

$$\Pr_{S, S'} [B(S, S') = 1] \leq$$
$$\underbrace{\Phi_d(2m) \cdot 2^{-\varepsilon m/2}}_{\substack{\text{\#labelings} \\ \text{of } T}} \Pr[B] \text{ wrt} \\ \text{permutation}$$

$$\leq C_0 (2m)^d 2^{-\varepsilon m/2}$$

(Sauer)

$$= e^{C_1 [d \log n - \varepsilon m]}$$

$$\Pr[A] \leq 2^x, \text{ set } \subseteq \mathcal{S} \text{ &} \\ \text{solve...}$$

Theorem Let $d = \text{VC}(\mathcal{H})$.

Then for

$$m \geq C_2 \left(\frac{1}{\varepsilon} \log \frac{1}{\delta} + \frac{d}{\varepsilon} \log \frac{1}{\varepsilon} \right)$$

with prob. $\geq 1 - \delta$

any consistent $h \in \mathcal{H}$

has $\Sigma(h) \leq \varepsilon$.

More generally,

if $\frac{1}{\varepsilon} \rightarrow \frac{1}{\varepsilon^2}$,

have uniform
convergence;

With prob. $\geq 1 - \delta$, $\forall h \in \mathcal{H}$:

$$|\hat{\varepsilon}_s(h) - \varepsilon(h)| \leq \varepsilon.$$

(slightly different $B(s, s')$
& prob. analysis)

Lower Bound #1

- Let $d = VC(\mathcal{H})$
- Let x_1, \dots, x_d be shattered
- $\forall \bar{v} \in \{0, 1\}^d$ let P_{uniform}
 - $h_{\bar{v}}(x_i) = v_i \quad \forall 1 \leq i \leq d$
- Now choose target $c = h_{\bar{v}}$
for \bar{v} random
- For $\epsilon \leq 1/4$, need $\Omega(d)$ examples - are just predicting coin flips!

Better lower
bound?

E.g. $\sqrt{d/\epsilon}$

for any ϵ ?

Lower Bound #2

- Same set-up as #1,
but now we'll
"give away" $c(x_i)$

P: $1 - \varepsilon$ $\overbrace{\quad \quad \quad}^{\varepsilon \text{ total, } \varepsilon/d-1 \text{ each}}$

• • • ... •
 x_1 x_2 x_3 x_d

- Still choose $c =$
random \bar{h}
- Now only see coin flip
every $\sim 1/\varepsilon$ samples \Rightarrow
 $\sqrt{n(d/\varepsilon)}$.

A Prescriptive

Application of VC:

Structural Risk

Minimization

- Suppose we have not a single \mathcal{H} but a nested hierarchy:

$$\mathcal{H}_1 \subset \mathcal{H}_2 \subset \mathcal{H}_3 \dots \subset \mathcal{H}_d \subset \dots$$

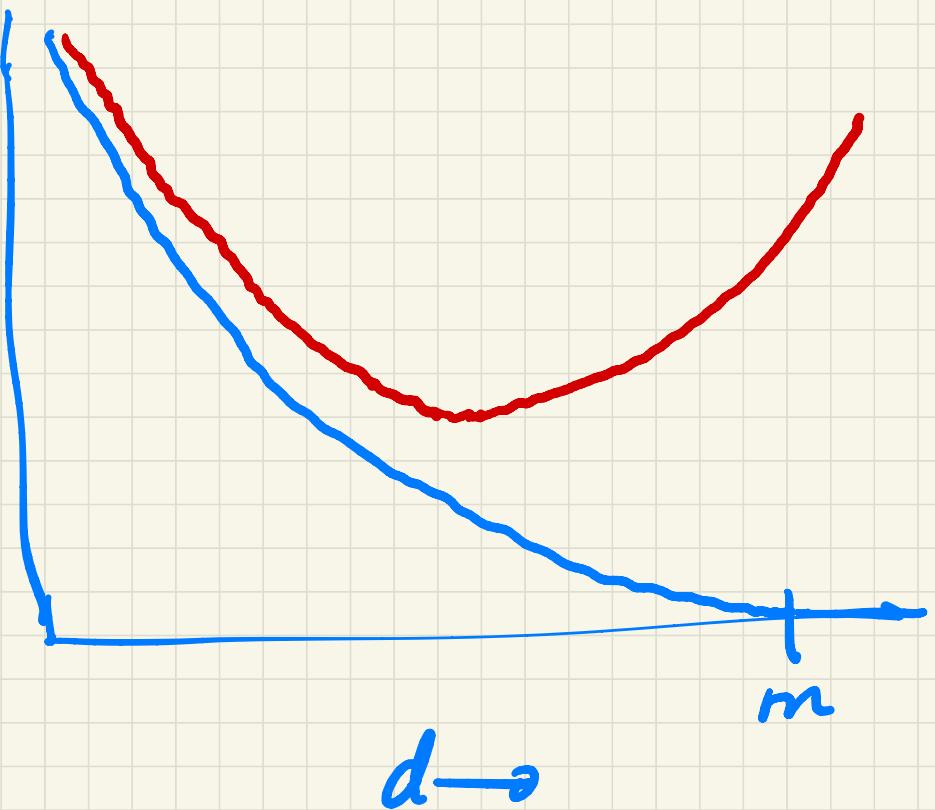
- E.g. neural networks of increasing depth/width
- For simplicity assume $VC(\mathcal{H}_d) = d$
- Let $S = \{(x_i, y_i)\}$ be of size m

VC theory says
that that in H_d ,

$$|\hat{\epsilon}_s(h) - \epsilon(h)| \leq \sqrt{\frac{d}{m}}$$

(ignoring log factors,
fixing δ , spreading
 δ over the H 's)

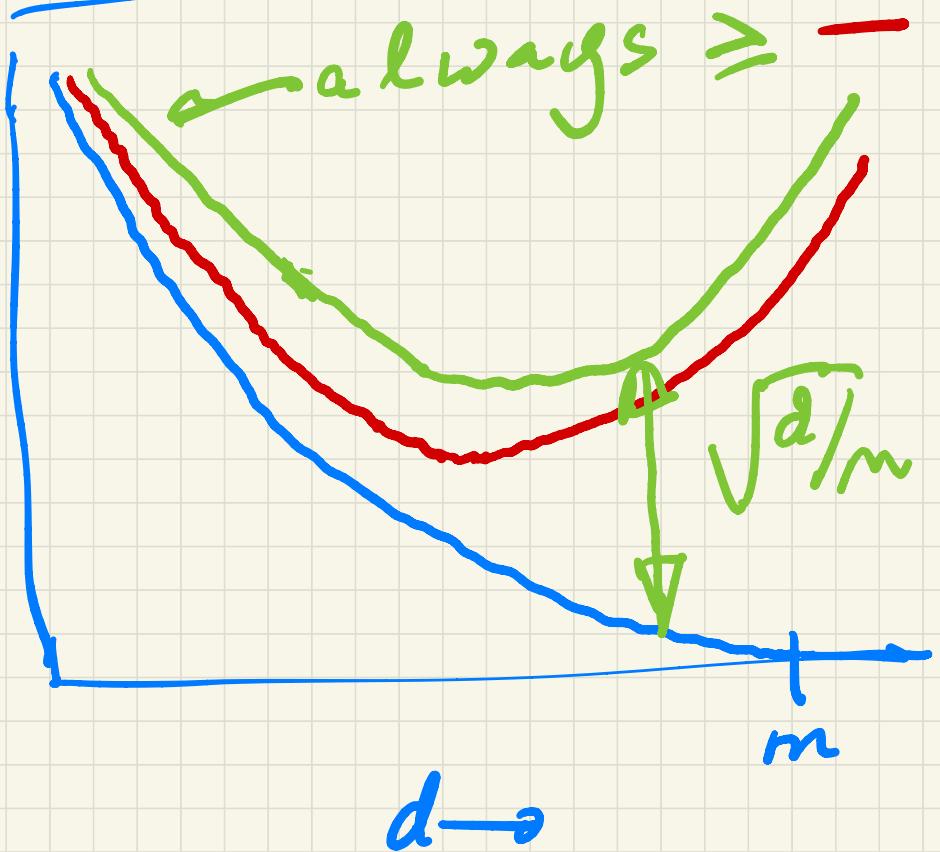
Classic overfitting cartoon:



—: optimal $\hat{\epsilon}_s(h)$ in \mathcal{H}_d

—: true $\epsilon(h)$ of

Classic overfitting cartoon:



—: optimal $\hat{\epsilon}_s(h)$ in H_d

—: true $\epsilon(h)$ of

—: $\hat{\epsilon}_s(h) + \sqrt{d/m}$

So choose:

$$d^* = \underset{d}{\operatorname{argmin}} \left\{ \hat{\epsilon}_S(h_d^*) + \sqrt{\frac{d}{m}} \right\}$$



$$h_d^* = \underset{h \in \mathcal{H}_d}{\operatorname{argmin}} \{ \hat{\epsilon}_S(h) \}$$

Then our true error

$$\leq \underset{d}{\operatorname{min}} \{ \text{same} \}$$

Distribution-Dependent
Improvements to
VC Theory

- Replace $\prod_{\mathcal{H}}(n)$ by
 $E_p[\prod_{\mathcal{H}}(s)]$ VC Entropy
- Related: Rademacher complexity
- Replace $|\mathcal{H}|(1-\varepsilon)^m$ by
 $\sum_{\varepsilon_i} |\mathcal{H}_i|(1-\varepsilon_i)^m$
 connections to
 stat mech

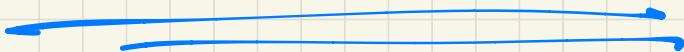
Data-Dependent
Improvements to
VC Theory

- Max-margin generalization
(e.g. linear, SVMs, boosting)

- "PAC-Bayes"

- small-norm generalization

Generalizations
to other
learning settings



- Models h in class \mathcal{H}
- Observations $z \sim P$ i.i.d.
- Loss function $\ell(h, z) \in \mathbb{R}$
also i.i.d

E.g. regression:

$$z = \langle x, y \rangle, y \neq h(x) \in \mathbb{R}$$

$$\ell(h, \langle x, y \rangle) = (h(x) - y)^2$$

(squared error)

$$\text{or } \ell(h, \langle x, y \rangle) = |h(x) - y|$$

(absolute err.)

$$\text{or } \ell(h, \langle x, y \rangle) =$$

$$(h(x) - y)^2 + \alpha |h(x)|$$

Distribution learning

- Models are distributions

$$P \in \mathcal{H}$$

- Observations $Z = X \cup P^*$
- Loss $\ell(P, x) = \log \frac{1}{P(x)}$

log-loss \Rightarrow MLE

Pretty much any model type, observations & loss function...

Now let $S = \{z_1, z_2, \dots, z_d\}$
and look at

$$\{\langle l(h, z_1), \dots, l(h, z_d) \rangle : h \in \mathcal{H}\}$$

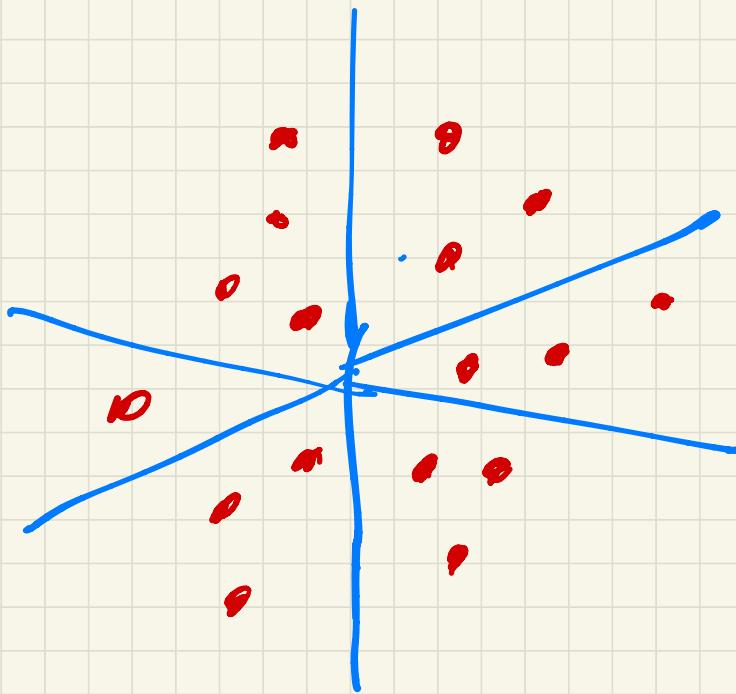
$\overbrace{\qquad\qquad\qquad}^{\in \mathbb{R}^d}$

= generalization of
 $\Pi_{\mathcal{H}}(S)$, now possibly
infinite!

Say \mathcal{H} shatters S if

$\Pi_{\mathcal{H}}(S)$ is
“space filling”...

E.g. $T_{\mathcal{H}}(s)$ intersects
all 2^d orthants of \mathbb{R}^d :



"fractal dimension"
"combinatorial dimension"

Moral: Uniform

convergence ↑
the norm, not

the exception.

