

12, 13.25 September 2021

CS 492(F), Comp. Learning Theory - Rademacher Complexity and VC-Dimension (Ch 3)

1. Motivation

(1) In the previous chapter, we learnt two results on uniform generalisation bounds for finite hypothesis sets, which let us state the guarantees of learning algorithms formally. However, these results stop applying as soon as the hypothesis set \mathcal{H} becomes infinite. In fact, this limitation cannot be fixed by naive tricks, because the conclusions of those results crucially rely on the finiteness of \mathcal{H} :

$$\mathbb{P} [\forall h \in \mathcal{H}, R(h) \leq \hat{R}_S(h) + \dots \log |\mathcal{H}| \dots + \dots] \geq 1 - \delta.$$

but $\log |\mathcal{H}| = \infty$ for an infinite \mathcal{H} .

(2) What should we do to handle infinite \mathcal{H} 's? The answer from this chapter is to use an appropriate notion of complexity of a hypothesis set and try to get a uniform generalisation bound using that notion, instead of $\log |\mathcal{H}|$.

(3) Rademacher complexity, growth function, and VC-dimensions are such notions. They all measure quantitatively the complexity of \mathcal{H} . We will study uniform generalisation bounds using these notions.

(4) The results of this chapter are not tied to specific learning algorithms. They are fundamental, and they are used repeatedly in the book to analyse concrete ML algorithms.

2. Setting

- (1) We assume general input and output spaces $X, Y \subseteq \mathbb{R}$.
Usually, Y is finite. For instance, $Y = \{-1, 1\}$ or $Y = \{0, 1\}$.
- (2) $\mathcal{H} \subseteq [X \rightarrow Y]$ is an assumed set of hypotheses.
- (3) $L : Y \times Y \rightarrow [0, 1]$ is an assumed loss function. Intuitively, it describes a loss. One example L is $(y_1, y_2) \mapsto \mathbb{1}_{\{y_1 \neq y_2\}}$.
- (4) $\mathcal{G} \stackrel{\text{def.}}{=} \{f : X \times Y \rightarrow \mathbb{R} \mid h \in \mathcal{H}\}$.

Note the close relationship between \mathcal{H} and \mathcal{G} . We will later show that the complexity of \mathcal{H} is closely related to that of \mathcal{G} .

3. Rademacher Complexity

- (1) \mathcal{U} ... input space $\mathcal{V} \subseteq \mathbb{R}$... output space.
 $\mathcal{F} \subseteq [\mathcal{U} \rightarrow \mathcal{V}]$... nonempty set of functions.

Def 3.1 Let $S \in \mathcal{U}^m$, that is, $S = (u_1, u_2, \dots, u_m)$.

The empirical Rademacher complexity of \mathcal{F} with respect to S is defined as follows:

$$\tilde{R}_S(\mathcal{F}) \stackrel{\text{def.}}{=} \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \xi_i f(x_i) \right] \quad \xi_i \sim \text{Uniform}\{-1, 1\}^m$$

Note that $\sum_{i=1}^m \xi_i f(x_i) = \langle (\xi_i)_{i \in [m]}, (f(x_i))_{i \in [m]} \rangle$ (i.e., inner product)

which gets maximized when $(f(x_i))_{i \in [m]}$ and $(\xi_i)_{i \in [m]}$ among $(f(x_i))_{i \in [m]}$'s costs $\| (f(x_i))_{i \in [m]} \| = c$ for some fixed c .

point to the same direction as m -dimensional vectors.

Thus, $\widehat{R}_S(F)$ measures the complexity of F by measuring how well F expresses different directions when we consider vectors $(f(x_i))_{i \in [m]}$ for $f \in F$.

E Def 3.2. Let $D \in \text{Pr}(\mathcal{U})$ and $m \in \mathbb{N}$. The Rademacher complexity of F with respect to D and m is

$$R_m(F) \stackrel{\text{def.}}{=} \mathbb{E}_{S \sim D^m} [\widehat{R}_S(F)],$$

that is, the expectation of the empirical Rademacher complexity over sample $S \sim D^m$.

(2) Example.

$$\textcircled{1} \quad F \subseteq \mathbb{R} \rightarrow \{1, -1\}$$

$$F = \{u \mapsto \mathbb{1}_{\{u \in [l, r]\}} - \mathbb{1}_{\{u \notin [l, r]\}} \mid l \leq r\}.$$

$$\textcircled{2} \quad S_1 = (3) \quad S_2 = (3, -2) \quad S_3 = (3, -2, 5).$$

$$\widehat{R}_{S_1}(F) = 1 \quad \widehat{R}_{S_2}(F) = 1 \quad \widehat{R}_{S_3}(F) = \frac{1}{8} \times 1 \times 7 + \frac{1}{8} \times \frac{2}{3} \times 1 = \frac{23}{24}$$

$$S_4 = (3, 3) \quad \widehat{R}_{S_4}(F) = \frac{1}{4} \times 1 \times 2 + \frac{1}{4} \times 0 \times 2 = \frac{1}{2}.$$

If D is the Dirac distribution on 3 (so that $(u_1, \dots, u_m) \sim D^m$ implies $u_1 = \dots = u_m = 3$), then $R_2(F) = \frac{1}{2}$. But if D is the uniform distribution on $[0, 1]$ (so that $(u_1, u_2) \sim D^2$ implies $u_1 \neq u_2$ with prob. 1), then $R_2(F) = 1$.

④ Note that R_m depends on the distribution D , and \widehat{R}_S depends on the sample S .

(3) Generalisation bound. Important boundedness requirement. Rademacher complexity.

Thm 3.3. $\mathcal{F} \subseteq \{u \rightarrow [0,1]\}$. That is, $D = [0,1]$.

Let $S > 0, m \in \mathbb{N}$ and D be a distribution on U . Then,

$$\underset{S \sim D^m}{\Pr} [\forall f \in \mathcal{F}, \mathbb{E}[f(u)] \leq \underset{\text{un } D}{\mathbb{E}}[f(u)] + 2R_m(\mathcal{F}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}]$$

$$\geq 1 - \delta.$$

↳ empirical distribution over S , i.e., uniform distribution over S .

$$\text{and } \underset{S \sim D^m}{\Pr} [\forall f \in \mathcal{F}, \mathbb{E}[f(u)] \leq \underset{\text{un } D}{\mathbb{E}}[f(u)] + 2\underset{\text{un } D_S}{R_S(\mathcal{F})} + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

↳ empirical Rademacher complexity.

which doesn't refer to a specific $f \in \mathcal{F}$

① $R_m(\mathcal{F})$ in the upper bound is what lets us have a uniform bound, i.e., a bound that works for all $f \in \mathcal{F}$.

② The proof uses McDiarmid's inequality.

Thm D.8 [McDiarmid's inequality]

$x_1, \dots, x_m \in \mathcal{X}$ for $m \geq 1$... independent random vars.

$c_1, \dots, c_m > 0$.

$$f: \mathcal{X}^m \rightarrow \mathbb{R} \text{ s.t. } |f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_m) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_m)| \leq c_i$$

for all x_1, \dots, x_m, x'_i, i .

$$\text{Then, } \Pr_{x_1, \dots, x_m} [f(x_1, \dots, x_m) - \mathbb{E}_{x_1, \dots, x_m}[f(x_1, \dots, x_m)] \geq \varepsilon] \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^m c_i^2}\right)$$

$$\text{and } \Pr_{x_1, \dots, x_m} [f(x_1, \dots, x_m) - \mathbb{E}_{x_1, \dots, x_m}[f(x_1, \dots, x_m)] \leq -\varepsilon] \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^m c_i^2}\right)$$

* More general than Hoeffding's inequality.

Proof of Thm 3.3. Fix $\delta > 0$, $m \in \mathbb{N}$.

Define $\Phi : \mathcal{U}^m \rightarrow [-1, 1]$ by

$$\Phi(S) \stackrel{\text{def.}}{=} \sup_{f \in \mathcal{F}} \left(\mathbb{E}_{\substack{u \sim D \\ u \in S}} [f(u)] - \mathbb{E}_{\substack{u \sim D \\ u \notin S}} [f(u)] \right).$$

for the first bound in the thm, it suffices to prove that

with prob. $\geq 1-\delta$ over $S \sim D^m$,

$$\Phi(S) \leq 2R_m(\mathcal{F}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

we will show: (i) $\mathbb{E}_{S \sim D^m} [\Phi(S)] \leq \mathbb{E}[\Phi(S)] + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$ w.p. $\geq 1-\delta$.

$$(ii) \mathbb{E}_{S \sim D^m} [\Phi(S)] \leq 2R_m(\mathcal{F}).$$

why does (i) hold? Because of McDiarmid's inequality as shown below:

$$\begin{aligned} & |\Phi(u_1, \dots, u_i, \dots, u_m) - \Phi(u_1, \dots, u'_i, \dots, u_m)| \\ & \leq \left| \sup_{f \in \mathcal{F}} \left(\mathbb{E}_{\substack{u \sim D \\ u \in \hat{D}_{(u_1, \dots, u_i, \dots, u_m)}}} [f(u)] - \mathbb{E}_{\substack{u \sim D \\ u \notin \hat{D}_{(u_1, \dots, u_i, \dots, u_m)}}} [f(u)] \right) \right| \\ & \quad - \left| \sup_{f \in \mathcal{F}} \left(\mathbb{E}_{\substack{u \sim D \\ u \in \hat{D}_{(u_1, \dots, u'_i, \dots, u_m)}}} [f(u)] - \mathbb{E}_{\substack{u \sim D \\ u \notin \hat{D}_{(u_1, \dots, u'_i, \dots, u_m)}}} [f(u)] \right) \right| \\ & = \left| \sup_{f \in \mathcal{F}} \frac{1}{m} (f(u'_i) - f(u_i)) \right| \leq \frac{1}{m} \end{aligned}$$

Thus, by McDiarmid's req., for any $\epsilon > 0$,

$$\mathbb{P}_{S \sim D^m} [\Phi(S) - \mathbb{E}_{S \sim D^m} [\Phi(S)] \leq \epsilon] \geq 1 - \exp \left(-\frac{2\epsilon^2}{m \cdot \left(\frac{1}{m} \right)^2} \right) = 1 - \exp(-2m\epsilon^2).$$

So, if $\epsilon = \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$, we can get (i).

called Rademacher variables.

What about (ii)? Because we can exploit symmetries of \vec{g} in the def'n of $R_m(\mathcal{F})$, as shown in the following calculation:

$$\begin{aligned}
\mathbb{E}_{S \sim D^m} [\Psi(S)] &= \mathbb{E}_{S \sim D^m} [\sup_{f \in \mathcal{F} \text{ w.r.t. } S} (\mathbb{E}_{u \sim D_S} [f(u)] - \mathbb{E}_{u \sim \hat{D}_S} [f(u)])] \\
&= \mathbb{E}_S [\sup_f (\mathbb{E}_{S' \sim D^m} [\mathbb{E}_{u \sim \hat{D}_{S'}} [f(u)]] - \mathbb{E}_{u \sim \hat{D}_S} [f(u)])] \\
&= \mathbb{E}_S [\sup_f \mathbb{E}_{S'} [\mathbb{E}_{u \sim \hat{D}_{S'}} [f(u)] - \mathbb{E}_{u \sim \hat{D}_S} [f(u)]]] \\
&\leq \mathbb{E}_{S, S'} [\sup_f (\mathbb{E}_{u \sim \hat{D}_{S'}} [f(u)] - \mathbb{E}_{u \sim \hat{D}_S} [f(u)])]. \\
&= \mathbb{E}_{S, S'} [\sup_f \frac{1}{m} \sum_{i=1}^m G_i (f(u_i) - f(u'_i))]. \\
&\stackrel{\text{G_i uniform $\mathbb{R}^{1,1^m}$}}{\leq} \mathbb{E}_{S, \vec{G}} [\sup_f \frac{1}{m} \sum_{i=1}^m G_i f(u'_i)] + \mathbb{E}_{S, \vec{G}} [\sup_f \frac{1}{m} \sum_{i=1}^m -G_i f(u_i)] \\
&\quad + \mathbb{E}_{S, \vec{G}} [\sup_f \frac{1}{m} \sum_{i=1}^m G_i f(u_i)] \\
&= 2 R_m(\mathcal{F}).
\end{aligned}$$

Flipping part
 of \$S\$ and \$S'\$
 based on \$G_i\$
 doesn't change
 the dist. of
 \$(S, S')\$

 \$(S, \vec{G})\$ and
 \$(S, -\vec{G})\$ have
 the same
 distributions.

the argument for.

For the second claim of the theorem, we repeat (i) but

for $\frac{\delta}{2}$, instead of δ . This gives us:

$$\mathbb{P}_{S \sim D^m} [\Psi(S) \leq \mathbb{E}_{S \sim D^m} [\Psi(S)] + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}] \geq 1 - \frac{\delta}{2}. \quad (\text{iii})$$

By (ii), we have $\mathbb{E}_{S \sim D^m} [\Psi(S)] \leq 2 R_m(\mathcal{F})$. (iv)

Now, if we apply McDiarmid's inequality to the function:

$$S = (u_1, \dots, u_m) \mapsto \widehat{R}_S(\mathcal{F})$$

then by essentially the same argument as the one for (i),

$$\mathbb{P}_{S \sim D^m} [R_m(\mathcal{F}) \leq \widehat{R}_S(\mathcal{F}) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}] \geq 1 - \frac{\delta}{2}. \quad (\text{v})$$

The claim of the theorem follows from (iii), (iv) and (v) via union bound. D.

4. Rademacher Complexity in Our Setup.

(i) How does the generalization bound apply in our setup?
in Thm 3.3

In our setup -

$$Y = \{-1, 1\}, \quad \mathcal{H} \subseteq \{\mathbb{X} \rightarrow Y\} \quad \text{.. hypo. set.}$$

$$\mathcal{U} = \mathbb{X} \times Y \quad \mathcal{F} = \{(\mathbf{x}, y) \mapsto \mathbb{1}_{\{h(\mathbf{x}) \neq y\}} \mid h \in \mathcal{H}\}.$$

Then, for $f(x, y) = \mathbb{1}_{\{h(x) \neq y\}}$ and $D \in \mathcal{P}_r(\mathcal{U}) = \Pr(\mathbb{X} \times Y)$

$$\mathbb{E}_{(x, y) \sim D} [f(x, y)] = \mathbb{E}_{(x, y) \sim D} [\mathbb{1}_{\{h(x) \neq y\}}] = R(h)$$

$$\mathbb{E}_{(x, y) \sim Q_S} [f(x, y)] = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x_i) \neq y_i\}} = \widehat{R}_S(h),$$

$(S = ((x_1, y_1), \dots, (x_m, y_m)))$

Thus, the bounds in Thm 3.3 become:

$$(vi) \underset{S \in \mathcal{D}^m}{\Pr} [\forall h \in \mathcal{H}. \quad R(h) \leq \widehat{R}_S(h) + 2R_m(\mathcal{F}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

$$(vii) \underset{S \in \mathcal{D}^m}{\Pr} [\forall h \in \mathcal{H}. \quad R(h) \leq \widehat{R}_S(h) + 2\widehat{R}_S(\mathcal{F}) + 3\sqrt{\frac{\log \frac{2}{\delta}}{2m}}] \geq 1 - \delta.$$

To get meaningful bounds, we should rephrase $R_m(\mathcal{F})$ and $\widehat{R}_S(\mathcal{F})$

in terms of $R_m(\mathcal{H})$ and $\widehat{R}_{S_{\mathbb{X}}}(\mathcal{H})$

*projection of S to \mathbb{X} -part.
That is, $S_{\mathbb{X}} = (x_1, \dots, x_m)$.*

(2)

Lemma 3.4 For all $S = ((x_1, y_1), \dots, (x_m, y_m)) \in (\mathbb{X} \times Y)^m$,

$$2\widehat{R}_S(\mathcal{F}) = \widehat{R}_{S_{\mathbb{X}}}(\mathcal{H})$$

where $S_{\mathbb{X}} = (x_1, \dots, x_m)$.

Proof $\widehat{R}_S(F) = \mathbb{E}_{\mathcal{S} \sim \text{Uniform } \{1..m\}^m} \left[\sup_{f \in F} \frac{1}{m} \sum_{i=1}^m G_i f(x_i, y_i) \right]$

(b_1, \dots, b_m) and ($y_1, b_1, \dots, y_m, b_m$) have the same distribution.

$$\begin{aligned}
 &= \mathbb{E}_{\mathcal{S}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m G_i \mathbb{I}_{\{h(x_i) \neq y_i\}} \right] \xrightarrow{\frac{1}{2} \text{ here is the reason that we can remove 2 from } 2\widehat{R}_S(F)} \\
 &= \mathbb{E}_{\mathcal{S}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m G_i \frac{1 - h(x_i) y_i}{2} \right] \\
 &= \mathbb{E}_{\mathcal{S}} \left[\frac{1}{m} \sum_{i=1}^m G_i \right] + \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m (-y_i b_i) h(x_i) \right] \\
 &= 0 + \frac{1}{2} \mathbb{E}_{\mathcal{S}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m G_i h(x_i) \right] = \frac{1}{2} \widehat{R}_{S\mathcal{H}}(\mathcal{H}) \quad \square
 \end{aligned}$$

Theorem 3.5 Let $D \in \mathcal{P}(X \times Y)$

Then, for all $\delta > 0$,

$$\Pr_{S \sim D^m} [\forall h \in \mathcal{H}, R(h) \leq \widehat{R}_S(h) + R_m(\mathcal{H}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta$$

and

$$\Pr_{S \sim D^m} [\forall h \in \mathcal{H}, R(h) \leq \widehat{R}_S(h) + \widehat{R}_{S\mathcal{H}}(\mathcal{H}) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}}] \geq 1 - \delta.$$

depend on D .

Proof. By Lemma 3.4, $2\widehat{R}_S(F) = \widehat{R}_{S\mathcal{H}}(\mathcal{H})$.

Taking the expectation over S gives: $2R_m(F) = R_m(\mathcal{H})$.

Replacing F -related parts in (vi) and (vii) using the above equations gives the conclusion of the theorem. \square .

5. Growth Function.

(1) Computing $R_S(\mathcal{H})$ and $R_m(\mathcal{H})$ is hard in most cases, partly because their dependence on sample S and distribution D . Actually, this of difficulty comes from the fact that they carry a lot of information about the problem. (which is a good thing.).

(2) The growth function and VC dimension are complexity measures of a hypothesis set \mathcal{H} , which are less informative than Rademacher complexity but easier to compute as a result. They describe how expressive \mathcal{H} is.

Def 3.6 The growth function $\Pi_{\mathcal{H}}: \mathbb{N} \rightarrow \mathbb{N}$ for a hypo. set.

\mathcal{H} is defined by

$$\Pi_{\mathcal{H}}(m) = \max_{\substack{(x_1, \dots, x_m) \in \mathcal{X}^m \\ h \in \mathcal{H}}} |\{h(x_1), \dots, h(x_m)\}|$$

① We continue to use the setup in 4. So, $\mathcal{Y} = \{-1, 1\}$.

② Note that $\Pi_{\mathcal{H}}(m) \leq 2^m$ (since $\mathcal{Y} = \{-1, 1\}$).

③ Large $\Pi_{\mathcal{H}}(m)$ means very expressive

Thm 3.7 [Massart's Lemma]. Let $A \subseteq \mathbb{R}^m$ be a finite nonempty

set. with $r = \max_{y \in A} \|y\|_2$. Then,

$$\mathbb{E}_{\substack{\text{uniform } \mathcal{A} \\ \tilde{y} \sim (\text{uniform } \{-1, 1\})^m}} \left[\frac{1}{m} \sup_{y \in A} \sum_{i=1}^m \tilde{y}_i y_i \right] \leq \frac{r \sqrt{2 \log |A|}}{m}$$

Pf. We use Cor D.11

Cor D.11 [Maximum inequality]

$(Y_{ij})_{(i,j) \in [m] \times [n]} \in \mathbb{R}^{m \times n}$... random variables s.t.

(i) $\mathbb{E}[Y_{ij}] = 0$, $Y_{ij} \in [r_{ij}, r_{ij}]$.

(ii) $Y_{1j}, Y_{2j}, \dots, Y_{nj}$ are independent
for every $j \in [n]$.

$$\text{Let } Z_j = \sum_{i=1}^m Y_{ij}.$$

Then,

$$\mathbb{E} \left[\max_{j \in [n]} Z_j \right] \leq r \sqrt{2 \log n}, \quad r = \sup_j \sqrt{\sum_{i=1}^m r_{ij}^2}$$

$(Y_{iy})_{(i,y) \in \mathbb{N}^m \times A}$, where $Y_{iy} = \epsilon_i y_i$ (random because $\epsilon_i \sim \text{Ber}(c)$)

Then Y_{iy} is mean zero, and $(Y_{iy})_i$ is indep. for each $y \in A$.

Also, $-|y_i| \leq Y_{iy} \leq |y_i|$. Thus, by Cor D.11,

$$\mathbb{E} \left[\max_{y \in A} \sum_{i=1}^m Y_{iy} \right] \leq \sup_{y \in A} \|y\|_2 \cdot \sqrt{2 \log |A|}$$

Thus

$$\mathbb{E} \left[\max_{y \in A} \frac{1}{m} \sum_{i=1}^m \epsilon_i y_i \right] \leq \frac{1}{m} \sup_{y \in A} \|y\|_2 \cdot \sqrt{2 \log |A|} . \quad \square$$

Cor 3.8. Let \mathcal{F} be a family of measurable functions

from \mathbb{X} to $\mathbb{S}_{1.13}$. Consider $D \in \mathcal{P}(\mathbb{X})$. Then,

$$R_m(\mathcal{F}) \leq \sqrt{\frac{2 \log \Pi_{\mathcal{F}}(m)}{m}}$$

Proof. For each $S = (x_1, \dots, x_m) \in \mathbb{X}^m$,

$$\mathcal{F}|_S \stackrel{\text{def}}{=} \{h(x_1), \dots, h(x_m) \mid h \in \mathcal{F}\}.$$

Then,

$$R_m(\mathcal{F}) = \mathbb{E} \left[\mathbb{E} \left[\sup_{h \in \mathcal{F}} \frac{1}{m} \sum \epsilon_i h(x_i) \right] \right]$$

$\sup_{S \in \mathbb{D}^m} \overline{\epsilon_i} \sim (\text{uniform } \mathbb{S}_{1.13})^m$

$$= \mathbb{E} \left[\mathbb{E} \left[\sup_{\substack{S \\ h \in \mathcal{F}|_S}} \frac{1}{m} \sum \epsilon_i y_i \right] \right] \leq \mathbb{E} \left[\frac{1}{m} \sqrt{m} \sqrt{2 \log |\mathcal{F}|_S|} \right]$$

$$\leq \sqrt{\mathbb{E} \left[\frac{2 \log \Pi_{\mathcal{F}}(m)}{m} \right]} = \sqrt{\frac{2 \log \Pi_{\mathcal{F}}(m)}{m}}$$

\square

Thm 3.9. For all $D \in \mathcal{P}(\mathbb{X} \times \mathbb{Y})$ and $S \in \mathbb{D}$,

$$\mathbb{P} \left[\forall h \in \mathcal{F}. R(h) \leq R_S(h) + \sqrt{\frac{2 \log \Pi_{\mathcal{F}}(m)}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

Prf. Immediate consequence of Thm 3.5 and Cor 3.8. R.

b. VC Dimension

- (1) The growth fn of \mathcal{H} is easier to deal with than the Rademacher complexity of \mathcal{H} , partly due to its independence on the distribution D on \mathcal{X} and sample $S \in \mathcal{X}^n$. But sometimes it is still difficult to compute. Can we get a measure on the complexity / expressiveness of \mathcal{H} that is easier to handle?
- (2) The VC dimension of \mathcal{H} is such a measure. It is more abstract (i.e., less informative) than the growth fn. $\Pi_{\mathcal{H}}$ by not being dependent on m , the size of sample. In so doing, it becomes a quantity that is easier to compute. Note the usual high-level strategy of using a less-informative but easier-to-handle concept

(3). $Y = \{-1, +1\}$. $\mathcal{H} \subseteq [\mathcal{X} \rightarrow Y]$
Def. A subset $X \subseteq \mathcal{X}$ is shattered if for every subset $X_0 \subseteq X$, there exists a hypo. $h \in \mathcal{H}$ s.t. $X_0 = \{x \in X \mid h(x) = +1\}$.

Def 3.10. The VC dimension of a hypothesis set \mathcal{H} is the size of the largest finite set that is shattered by \mathcal{H} :
$$\text{VCdim}(\mathcal{H}) = \sup \{m \mid \Pi_{\mathcal{H}}(m) = 2^m\}$$

① Note that the largest finite set might not exist. To cover such a case, we use \sup here. In that case, $\text{VCdim}(\mathcal{H}) = \infty$.

② $\text{VC Dim } (\mathcal{H}) = d$ doesn't imply that every subset of \mathcal{X} with size d is shattered. Instead, it means some subset is shattered.

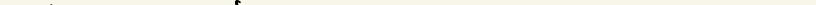
③ $\text{VCDim}(\mathcal{H}) = d$ implies that for every $d' \leq d$, there is an \mathcal{H} -shattered subset of \mathcal{X} with size d' . Exercise: prove this.

④ If \mathcal{H} consists of diverse fns, for a given set of labelled inputs $(x_1, y_1), \dots, (x_m, y_m)$, \mathcal{H} is likely to have a hypothesis h that has zero error for this set. The VC dimension uses this informal consequence to measure the complexity of \mathcal{H} .
 $VC(\mathcal{H}) = d$ means that for some size- d subset $X_0 \subseteq \mathcal{X}$, however we label inputs in X_0 , we can achieve 0 error with some $h \in \mathcal{H}$.

① Intervals on the real line.

$$X = \mathbb{R} \quad , \quad f(x) = \begin{cases} x & \text{if } l \leq x \leq u \text{ then } 1 \\ -1 & \text{else} \end{cases} .$$

$$\text{VC Dim}(\text{ff}) = 2 \quad \text{why?} \quad \begin{array}{c} \bullet \quad \bullet \\ \boxed{+} \quad \boxed{+} \quad \boxed{-} \\ \boxed{+} \quad \boxed{-} \quad \boxed{+} \\ \boxed{-} \quad \boxed{+} \end{array} \quad \mathbb{R} \quad \text{So, shattered.}$$

But:  cannot be R

classified correctly by any $h \in \mathcal{H}$.

② Hyperplanes in \mathbb{R}^d . $\text{Sgn}(y) = \begin{cases} +1 & \text{if } y \geq 0 \\ -1 & \text{if } y < 0 \end{cases}$

$$x \in \mathbb{R}^d, \quad f(x) = \begin{cases} 1 & \text{if } x \in \{x \mid \omega \cdot x + b \geq 0\} \\ -1 & \text{if } x \in \{x \mid \omega \cdot x + b < 0\} \end{cases}$$

$$\text{VCDim}(\mathcal{F}) = d+1$$

(i) why $\text{VCDim}(\mathcal{H}) \geq d+1$? $x_0 = (0, 0, \dots, 0)$
 $x_1 = (1, 0, \dots, 0)$

$$x_0 = (0, 0, \dots, 0) \in \mathbb{R}^d$$

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

$$x_2 = (0, 0, \dots, 0, -1)$$

Then, for every $(y_0, \dots, y_d) \in \{-1, +1\}^{d+1}$, pick the following $n \in \mathbb{N}$:

$$h(x) = \text{sgn}(\langle y, x \rangle + \frac{y_0}{2}) \quad \text{where } y = (y_1, \dots, y_d).$$

Then, $h(x_0) = \text{sgn}(\frac{y_0}{2}) = y_0$. and $h(x_i) = \text{sgn}(y_i + \frac{y_0}{2}) = y_i$ ($i \geq 1$).

(ii) why $\text{VCdim}(h) \leq d+1$? Because of Radon's thm.

Thm 3.13 [Radon's thm] Any $X \subseteq \mathbb{R}^d$ with $|X| = d+2$.

can be partitioned into X_1 and X_2 s.t. the convex hulls of X_1 and X_2 overlap.

input \downarrow

every input \downarrow

Note that if we label every x_i with +1 and x_2 with -1,

the X with this labelling cannot be classified correctly by any $h \in H$.

Thus, the thm implies that $\text{VCdim}(h) \leq d+2$.

Proof. Let $X = \{x_1, \dots, x_{d+2}\} \subseteq \mathbb{R}^d$ and consider

the following system of $(d+1)$ linear equations over $d+2$ variables

$$d_1, \dots, d_{d+2} : \begin{array}{c} d \text{ equations} \\ \sum_{i=1}^{d+2} d_i x_i = \vec{0} \quad (\in \mathbb{R}^d) \end{array} \quad \sum_{i=1}^{d+2} d_i = 0 \quad \dots \quad 1 \text{ equation}$$

Then, \exists a solution $(\beta_i)_{i \in [d+2]}$ s.t. $\beta_i \neq 0$ for some $i \in [d+2]$.

Let $J_1 = \{j \mid \beta_j > 0\}$ and $J_2 = \{j \mid \beta_j \leq 0\}$.

Since $\sum_{j=1}^{d+2} \beta_j = 0$ and $\beta_i \neq 0$ for some i , $J_1 \neq \emptyset$ and $J_2 \neq \emptyset$.

Let $X_1 = \{x_j \mid j \in J_1\}$ and $X_2 = \{x_j \mid j \in J_2\}$.

Define $\beta = \sum_{j \in J_1} \beta_j$. Then, since $\sum_{i=1}^{d+2} \beta_i x_i = \vec{0}$,

$$\sum_{j \in J_1, \beta} \frac{\beta_j}{\beta} x_j = \sum_{j \in J_2} -\frac{\beta_j}{\beta} x_j.$$

This means that $y = \sum_{j \in J_1, \beta} \frac{\beta_j}{\beta} x_j$ belongs to the convex hull of X_1 and also to that of X_2 .

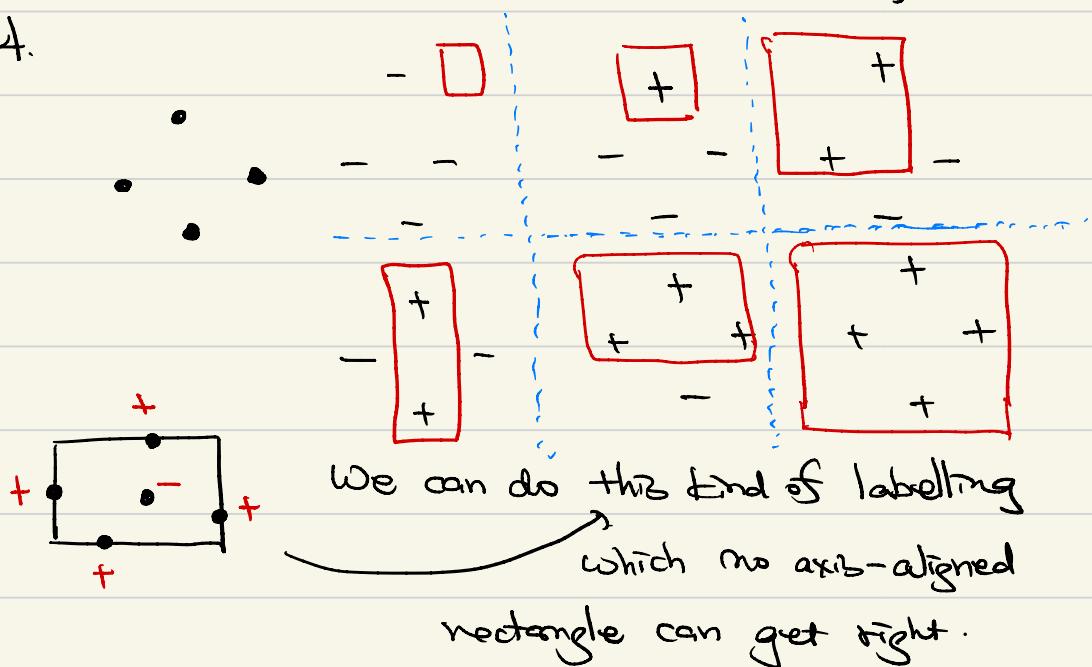
Q.E.D.

③ Axis-aligned rectangle.

$$X = \mathbb{R}^2, \quad f(x) = \begin{cases} 1 & \text{if } x \in [l-r] \times [b, t] \\ -1 & \text{else} \end{cases}$$

$$\text{vCDim } (\text{ff}) = 4.$$

Why ≥ 4 ?



④ Convex Polygon with d vertices.

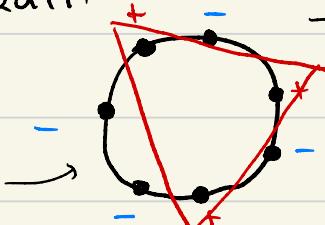
$$\text{VCDim}(\mathcal{F}) = 2d+1.$$

Why $\Sigma 2d+1$?

equally spaced. Points

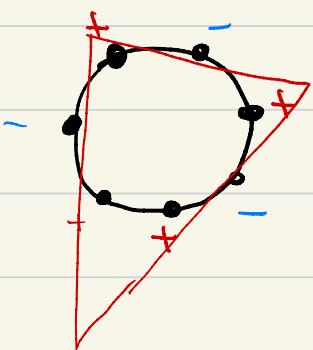
$\mathcal{F} = \mathbb{R}^2$ \mathcal{F} = fns generated
by convex d-gons.
fewer positive points.

Start from a triangle touching positive points and extend it.



fewer negative parts.

Start from a triangle
touching negative points and
shift it



⑤ Some functions. $\mathbb{X} = \mathbb{R}$, $f = \{x \mapsto -\operatorname{sgn}(\sin(\omega x)) \mid \omega \in \mathbb{R}\}$.

$$VC\dim(\mathcal{H}) = \infty.$$

(5) Generalization bound with VC dim.

① Assume that $\text{VC Dim}(\text{ff}) = d$. Then,

$$\Pi_{\text{fp}}(m) \leq \left(\frac{em}{d}\right)^d \quad (\text{Cor. 3.18}) \quad \text{VC-dim.}$$

Thus, we can derive a uniform generalization bound involving d .

from the one based on the growth fn., which we recall below.

$$\Pr_{S \sim D^m} [\forall h \in \text{ff. } R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2 \log \Pi_{\text{fp}}(m)}{m}} + \sqrt{\frac{\log 1/\delta}{m}}] \geq 1 - \delta.$$

where D is a dist. over $\mathcal{X} \times \mathcal{Y}$. That is, we can get:

$$\Pr_{S \sim D^m} [\forall h \in \text{ff. } R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2d \log \frac{em}{d}}{m}} + \sqrt{\frac{\log 1/\delta}{m}}] \geq 1 - \delta.$$

(Cor 3.19)

② There is more detailed and more formal explanation.

Thm 3.17 [Sauer's Lemma].

$\text{VC Dim}(\text{ff}) = d$ implies that $\Pi_{\text{fp}}(m) \leq \sum_{i=0}^d \binom{m}{i}$ for all m .

Proof. Induction on $m+d$.

Base case: $d=0$. Then, $|\text{ff}|=1$. Thus, $\Pi_{\text{fp}}(m)=1$, and the thm holds.
 $m=1$ and $d>0$. Then, $\sum_{i=0}^d \binom{m}{i} = \sum_{i=0}^d \binom{1}{i} \geq 2$. But $\Pi_{\text{fp}}(1) \leq 2$.

So, the thm holds in this case as well.

Inductive Case: $d \geq 0$ and $m \geq 1$.

for $S_0 \subseteq \mathcal{X}$, $\text{ff} \Big|_{S_0} \stackrel{\text{def.}}{=} \{h \mid S_0 \in [S_0 \rightarrow \mathcal{Y}] \mid h \in \text{ff}\}$.

Pick $S \subseteq \mathcal{X}$ with $|S|=m$ s.t. $\Pi_{\text{fp}}(m) = |\text{ff}|_S|$.

Let (x_1, \dots, x_m) be the enumeration of elements of S .

Let $S' = \{x_1, \dots, x_{m-1}\}$. (i.e., $S' = S \setminus \{x_m\}$).

and $\mathcal{H}' = \{h \in \mathcal{H}|_{S'} \mid \exists h_1, h_2 \in \mathcal{H}|_S \text{ s.t.}$

$$h = h_1|_{S'} = h_2|_{S'} \text{ and}$$

$$h_1(x_m) = +1 \text{ and } h_2(x_m) = -1\}$$

Then,

$$|\mathcal{H}'| + |\mathcal{H}|_{S'}| = |\mathcal{H}|_S| = \Pi_{\mathcal{H}}(m).$$

$$\text{Also, } \text{VCDim}(\mathcal{H}|_{S'}) \leq \text{VCDim}(\mathcal{H}) = d \text{ and}$$

$$\text{VCDim}(\mathcal{H}') \leq \text{VCDim}(\mathcal{H}|_S) - 1 = \text{VCDim}(\mathcal{H}) - 1 = d - 1.$$

Thus,

$$\Pi_{\mathcal{H}}(m) = |\mathcal{H}'| + |\mathcal{H}|_{S'}|$$

$$= \Pi_{\mathcal{H}'}(m-1) + \Pi_{\mathcal{H}|_{S'}}(m-1)$$

$$\begin{aligned} \text{Ind. hypo.} \rightarrow & \leq \sum_{i=0}^{d-1} \binom{m-1}{i} + \sum_{i=0}^d \binom{m-1}{i} = \sum_{i=0}^d \binom{m-1}{i-1} + \binom{m-1}{i} \\ & = \sum_{i=0}^d \binom{m}{i} \end{aligned}$$

D.

Cor 3.18. $\text{VCDim}(\mathcal{H}) = d$ implies that

$$\Pi_{\mathcal{H}}(m) = \left(\frac{e^m}{d}\right)^d = O(m^d)$$

for all $m \geq d$.

$$\stackrel{\text{Pf.}}{\Rightarrow} \Pi_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i}$$

$$\leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i}$$

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

$$\leq \left(\frac{m}{d}\right)^d e^{m \cdot \frac{d}{m}} = \left(\frac{e^m}{d}\right)^d$$

D.

Cor 3.19.

$$\text{VC Dim (ff)} = d \quad . \quad S > 0 \quad . \quad D \in \Pr(\mathcal{X} \times \mathcal{Y}).$$

$$\Rightarrow \Pr_{\mathcal{S} \sim D^m} [\text{the ff. } R(h) \leq \hat{R}_S(h) + \sqrt{\frac{d \log(\epsilon m/d)}{m}} + \sqrt{\frac{\log \epsilon}{2m}}] \geq 1 - \delta.$$

Proof

By Cor. 3.18 and Cor 3.9.

D.

$$R(h) \leq \hat{R}_S(h) + O\left(\sqrt{\frac{\log m/d}{m/d}}\right).$$

So the smaller d , the tighter the bound. Related to Occam's razor principle.