

# LECTURE 25

## SUMMARY of math framework of ML <sup>(Last class)</sup>

- $\mathcal{P}$ : an unknown distribution on  $X \times Y$ ; ↙ sets
- We see training data:  $(x_1, y_1), \dots, (x_n, y_n) \sim \mathcal{P}$  iid. ↙ label Goal = oracle  $h: X \rightarrow Y$ :  $h(x) = y$
- Choose a hypothesis class  $\mathcal{H}$  (functions  $X \rightarrow Y$ )
- $\forall h \in \mathcal{H}$ , Risk, a.k.a. "test error":  
↙ loss function, e.g.  $\mathbb{E}(h(x) - y)^2$  (quadratic loss)  
 $R(h) := \mathbb{E} \ell(h(x), y)$        $h^* = \operatorname{argmin}_{h \in \mathcal{H}} R(h)$ .      Not computable
- Empirical risk a.k.a. training error
- $R_n(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$ .       $h_n^* := \operatorname{argmin}_{h \in \mathcal{H}} R_n(h)$ .      Computable
- ERM algorithm:  

- ① Training: for input data  $(x_1, y_1), \dots, (x_n, y_n)$ ; compute  $h_n^*$ .
  - ② Prediction: on query  $X$ , output  $h_n^*(X)$  "oracle"

Lemma (Generalization error)

$$R(h_n^*) \leq R(h^*) + 2 \sup_{h \in \mathcal{H}} |R_n(h) - R(h)|$$

↑ test error of ERM      ↑ best possible error (with  $\infty$  data) for a given class  $\mathcal{H}$

Proof

$$\begin{aligned}
 R(h_n^*) &\leq R_n(h_n^*) + \epsilon && (h_n^* \in \mathcal{H}) \\
 &\leq R_n(h^*) + \epsilon && (h^* = \text{minimizer of } R_n) \\
 &\leq R(h^*) + 2\epsilon && (h^* \in \mathcal{H}) \quad \square.
 \end{aligned}$$

Next time, add a lemma:

if  $\|f - g\|_\infty < \epsilon$   
 and  $x^*, y^*$  are  
 minimizers of  $f, g$ ,  
 then  
 $|f(x^*) - g(y^*)| < 2\epsilon$ .

Then apply this lemma  
 for  $f = R_n, g = R$ ,  
 $h_n^*, h^*$

$$\epsilon = \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n Z_i(h) - \mathbb{E}Z(h) \right| \quad \text{where } Z_i(h) = \ell(h(x_i), y_i)$$

are iid r.v.'s  $\forall h$ .

"empirical process"  $\parallel$   
 $\bar{Z}_i(h)$

• For binary classification,  $\ell(\cdot, \cdot) \in \{0, 1\} \Rightarrow |Z_i(h)| \leq 1$

$$P\left\{ \left| \frac{1}{n} \sum_{i=1}^n \bar{Z}_i(h) \right| > t \right\} = P\left\{ \left| \frac{1}{n} \sum_{i=1}^n \bar{Z}_i(h) \right| > t/n \right\} \leq 2 \exp\left(-\frac{(t/n)^2}{2}\right) \quad (*)$$

multiply both sides by  $\sqrt{n}$  to scale like in CLT iid mean 0, bdd by 1 General Hoeffding inequality (lec.5)

• Union Bound:

$$P\left\{ \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \bar{Z}_i(h) \right| > t \right\} \leq \sum_{h \in \mathcal{H}} P\left\{ \left| \frac{1}{n} \sum_{i=1}^n \bar{Z}_i(h) \right| > t \right\}$$

$\leftarrow$  assume  $\mathcal{H}$  is finite  
 $\leftarrow$  "exists  $h \in \mathcal{H}$ " = union

$$\stackrel{(*)}{\leq} |\mathcal{H}| \cdot 2 \exp(-t^2/n/2) = 2 \exp(\log|\mathcal{H}| - t^2/n/2)$$

$$\Rightarrow \text{We proved:} \quad \leq 0.01 \quad \text{if } t = C \sqrt{\frac{\log|\mathcal{H}|}{n}}$$

THM (Generalization Bound) If the hypothesis class  $\mathcal{H}$  is finite,

$$R(h_n^*) \leq R(h^*) + C \sqrt{\frac{\log|\mathcal{H}|}{n}} \quad \text{with prob. } \geq 0.99.$$

$\uparrow$  ERM's test error  $\uparrow$  best possible error

$$\underbrace{\hspace{10em}}_{0.01} \quad \text{if } n \geq C' \log|\mathcal{H}|$$

Hence the ERM algorithm generalizes well from

$$n \sim \log|\mathcal{H}|$$

training data points.

• Good: logarithmic in  $|\mathcal{H}|$

• Bad: most hypothesis classes are infinite.

Can  $\log|\mathcal{H}|$  be replaced by some "complexity" of  $\mathcal{H}$ ?

Yes: VC dimension.

# VC DIMENSION

(Ванник - Черволенкис)

• Heuristically:  $vc(\mathcal{H}) = \text{largest } \#(\text{data } \mathcal{H} \text{ overfits})$

↖ i.e. functions  $h: X \rightarrow \{0,1\}$

**Def** Let  $\mathcal{H}$  be any collection of Boolean functions on a set  $X$ .

We say that  $\mathcal{H}$  overfits, or "shatters" a subset  $\{x_1, \dots, x_d\} \subset X$

if  $\forall$  labels  $y_1, \dots, y_d \in \{0,1\} \exists h \in \mathcal{H}$  such that

$$h(x_i) = y_i \quad \forall i=1, \dots, d.$$

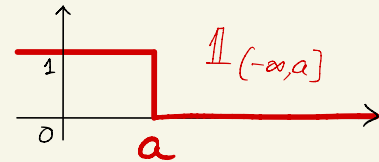
The **vc dimension** of  $\mathcal{H}$ , denoted  $vc(\mathcal{H})$ , is the maximal size  $d$  of a subset  $\mathcal{H}$  shatters.

## Examples

↖  $h(x) = 1 \quad \forall x$

1.  $\mathcal{H} = \{ \mathbb{1} \}$  has  $vc(\mathcal{H}) = 0$ : it can't shatter even one point  $x_i$  since  $h(x_i) = 1$ .

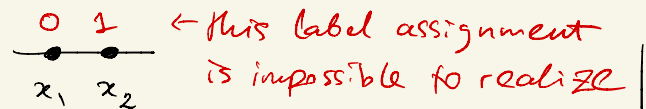
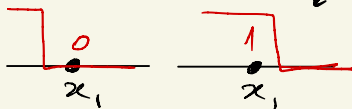
2. Half-lines  $\mathcal{H} = \{ \mathbb{1}_{(-\infty, a]} : a \in \mathbb{R} \}$



$$vc(\mathcal{H}) = 1$$

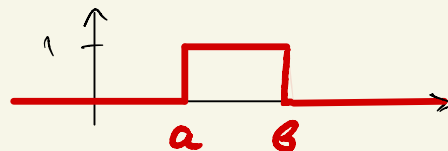
HW:  $\{ \mathbb{1}_{(-\infty, a]} ; \mathbb{1}_{[b, +\infty)} \}$

**Proof:**  $\mathcal{H}$  can shatter some 1-point set  $\{x_1\}$ , but can't shatter any 2-point set  $\{x_1, x_2\}$

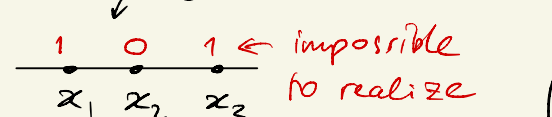
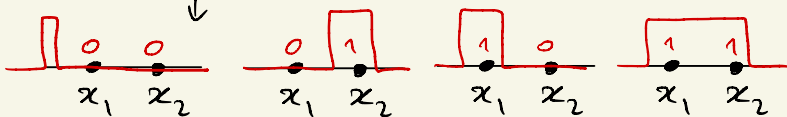


3. Intervals:  $\mathcal{H} = \{ \mathbb{1}_{[a, b]} : a \leq b \}$

$$vc(\mathcal{H}) = 2$$

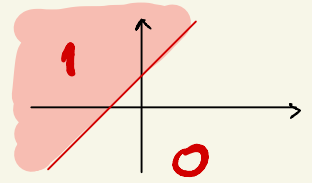


**Proof:**  $\mathcal{H}$  can shatter some 2-pt set  $\{x_1, x_2\}$ , but can't shatter any 3-point set  $\{x_1, x_2, x_3\}$



#### 4. Half-planes in $\mathbb{R}^2$ :

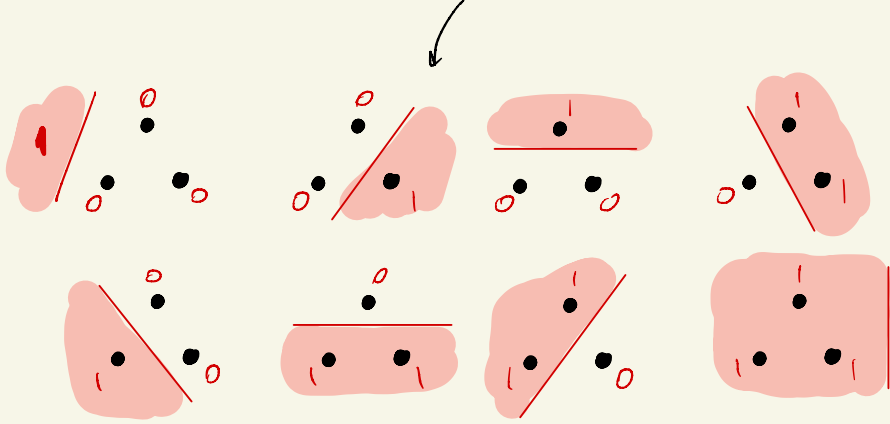
$$\mathcal{H} = \left\{ \mathbb{1}_{\{a_1 x(1) + a_2 x(2) + b \geq 0\}} : a_1, a_2, b \in \mathbb{R} \right\}$$



$$vc(\mathcal{H}) = 3$$

<sup>Proof</sup>  
 $\mathcal{H}$  can shatter some 3-pt set  $\{x_1, x_2, x_3\}$ ,

but can't shatter any 4-point set  $\{x_1, x_2, x_3, x_4\}$ :



A 4-point set is  
 like this, or like this  
  
 "Convex position"      "non-convex position"

In either case,  $\exists$  label assignment that is impossible to realize