

# Uniform convergence and Rademacher Complexity

## I) Introduction

- Reminder on Hoeffding inequality

$$P\left(\left|\frac{1}{N} \sum Z_i - \mathbb{E}Z\right| \geq \varepsilon\right) \leq 2 \exp(-2N\varepsilon^2)$$

equivalently,

$$\left|\frac{1}{N} \sum Z_i - \mathbb{E}Z\right| < \sqrt{\frac{\log 2/\delta}{2N}} \quad \text{w.p. at least } 1-\delta$$

- def: a sequence of r.v.  $Z_1, \dots, Z_N, \dots$  converges in proba to  $Z$  ( $Z_N \xrightarrow[N \rightarrow \infty]{\text{proba}} Z$ ) iff:

$$\forall \varepsilon, \delta \in ]0, 1[, \exists n, \text{ if } N > n,$$

$$|Z_N - Z| < \varepsilon \text{ with proba } 1-\delta$$

- Equivalently, there exist a function  $n(\varepsilon, \delta)$  such that  $\forall \varepsilon, \delta \in ]0, 1[,$

$$N > n(\varepsilon, \delta) \Rightarrow |Z_N - Z| < \varepsilon \text{ with proba } 1-\delta$$

exercise: Show in the Hoeffding setting

$$\frac{1}{N} \sum_i Z_i \xrightarrow[N \rightarrow \infty]{\text{in proba}} \mathbb{E}Z, \text{ give } n(\varepsilon, \delta)$$

- In this lecture,  $S = \{(x_1, y_1), \dots, (x_N, y_N)\}$  is considered as a random variable

- The empirical risk

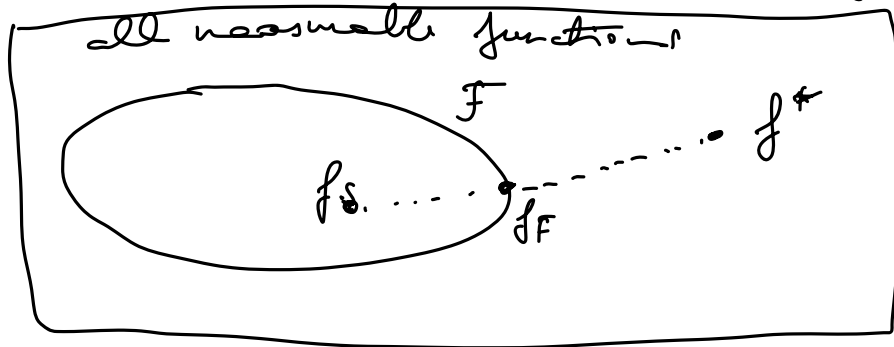
$$\hat{R}_S(f_S) = \frac{1}{N} \sum_{i=1}^N \ell(f_S(x_i), y_i)$$

↑ can be the 0/1 loss

## II Notions of consistency

A learner  $f_S$  is  $ERN$  iff

$$f_S \in \arg \min_{f \in F} \hat{R}_S(f)$$



$$f^* = \arg \min_{f \in \text{measurable}} R(f)$$

$$f_F \in \arg \min_{f \in F} R(f)$$

$$R(f_S) \geq R(f_F) \geq R(f^*)$$

estimation error
approximation error

Def: The learning algo  $f_S$

- is universally Bayes consistent iff  
 $\forall \text{ distribution } \mathcal{P}, \quad R(f_S) \xrightarrow[N \rightarrow \infty]{\text{in prob.}} R(f^*)$

in other words, there is a function  $m(\epsilon, \delta, \mathcal{P})$  such that for any  $\mathcal{P}, \epsilon, \delta$   
 if  $N > m(\epsilon, \delta, \mathcal{P})$  then,  
 for  $S \sim \mathcal{P}^N$ ,  $|R(f_S) - R(f^*)| < \epsilon$  w.p.  $1 - \delta$

$\triangle$  impossible for  $ERN$

- is universally  $F$ -consistent if  
 $\forall \mathcal{P}, \quad R(f_S) \xrightarrow[N \rightarrow \infty]{\text{in pr.}} R(f_F)$
- is a PAC-learner (Probably approximately correct)  
 if there is a function  $m(\epsilon, \delta)$

such that  $\forall$  distribution  $P$   
 $\forall \epsilon, \delta \in ]0, 1[$  - if  $N > m(\epsilon, \delta)$   
 then for  $S \sim P^N$ ,  $|R(f_S) - R(f_F)| < \epsilon$  w.p.  $1 - \delta$   
 $\triangle$  PAC implies F consistency

### III) PAC learning and Uniform convergence for ERM

- We want to bound  $R(f_S) - R(f_F)$
- Hoeffding allows us to bound  $R(f) - \hat{R}(f)$   
for a fixed  $f$  !!!

$$\begin{aligned} R(f_S) - R(f_F) &= R(f_S) - \hat{R}(f_S) + \underbrace{\hat{R}(f_S) - \hat{R}(f_F)}_{\leq 0} + \hat{R}(f_F) - R(f_F) \\ &\leq 2 \sup_{f \in F} |R(f) - \hat{R}(f)| \end{aligned}$$

def: The Unrepresentativeness of  $S$  w.r.t.  $F$  is  

$$\text{Unrep}(F, S) = \sup_{f \in F} |R(f) - \hat{R}(f)|$$

thm: If, for class  $F$ , there exist  $m(\epsilon, \delta)$  s.t.  
 for any distribution  $P$ , any  $\epsilon, \delta \in ]0, 1[$ ,  
 if  $N > m(\epsilon, \delta)$  then  $\text{Unrep}(F, S) < \epsilon$  w.p.  $1 - \delta$   
 (which is called the uniform convergence property),  
 then, ERM is a PAC learner on  $F$ .

proof: if  $N > m(\frac{\epsilon}{2}, \delta)$  then  $\text{Unrep}(F, S) \leq \frac{\epsilon}{2}$  w.p.  $1 - \delta$   
 and  $R(f_S) - R(f_F) \leq 2 \text{Unrep}(F, S) \leq \epsilon$  w.p.  $1 - \delta$ , so  $f_S$  is a PAC learner

## Application to finite class $F$

I want to show  $\sup_{f \in F} |R(f) - \hat{R}(f)| < \epsilon$  w.p.  $1 - \delta$   
for  $N > n(\epsilon, \delta)$ .

$$\begin{aligned} & P\left(\sup_{f \in F} |R(f) - \hat{R}(f)| \geq \epsilon\right) \\ &= P(\exists f \in F, |R(f) - \hat{R}(f)| \geq \epsilon) \end{aligned}$$

$$\leq \sum_{f \in F} P(|R(f) - \hat{R}(f)| \geq \epsilon)$$

*this  $f$  does not depend on the data  
 $\Rightarrow$  Hoeffding!*

$$\text{note that: } \hat{R}(f) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i), y_i)$$

$$R(f) = \mathbb{E}_{S \sim P^n} [\hat{R}(f)]$$

$$P\left(\sup_{f \in F} |R(f) - \hat{R}(f)| \geq \epsilon\right) \leq \sum_{f \in F} \overbrace{P(|R(f) - \hat{R}(f)| \geq \epsilon)}^{\text{apply Hoeffding here}} \leq \delta \times |F|$$

if  $N > \frac{\log \frac{2}{\delta}}{2\epsilon^2}$

$$\Rightarrow \text{Unif}(F, S) \leq \epsilon \text{ w.p. } 1 - \delta \times |F| \text{ when } N > \frac{\log 1/\delta}{2\epsilon^2}$$

$$\text{let } \delta' = \delta \times |F|$$

$$\Rightarrow \text{Unif}(F, S) \leq \epsilon \text{ w.p. } 1 - \delta' \text{ when } N > \frac{\log \left(\frac{2|F|}{\delta'}\right)}{2\epsilon^2}$$

$\Rightarrow$  ERM on finite class is PAC-learnable.

Equivalently,

$$\text{Unrep}(F, S) \leq \sqrt{\frac{\log 2|F|}{2N}} \quad \text{w.p. at least } 1-\delta'$$

$$\Rightarrow |R(f_S) - \hat{R}(f_S)| \leq \sqrt{\frac{\log 2|F|}{N}} \quad \text{w.p. at least } 1-\delta'$$

⚠ This only works for finite class  
because the union bound for  $\infty$  nb  
of events looks like  $P(\exists i, A_i) \leq \sum_{i=1}^{\infty} P(A_i) \approx \infty$

#### IV The case $|F| = \infty$ , Rademacher complexity

Goal: bound  $\text{Unrep}(F, S)$  for  $|F| = \infty$   
without union bound

Many tools: Vapnik dim, covering numbers,  
gaussian complexity, Rademacher  
complexity, ...

Rademacher applies to arbitrary bounded losses.

Notation:  $Z = (X, y)$  a labelled example  
 $S = (Z_1, \dots, Z_N)$

Given  $F$

define  $G = \ell \circ F = \{(x, y) \mapsto \ell(f(x), y), f \in F\}$

$$\begin{aligned} \text{UnRep}(F, S) &= \sup_{f \in F} |R(f) - \hat{R}(f)| \\ &= \sup_{g \in G} \left| \frac{1}{N} \sum_{i=1}^N g(Z_i) - \mathbb{E}_{Z \sim P} g(Z) \right| \end{aligned}$$

def: the empirical Rademacher complexity of  $S$  on  $g$  is

$$\hat{Rad}_S(G) = \frac{1}{N} \mathbb{E}_{\sigma_1 \dots \sigma_N \sim \text{unif}(\{-1, 1\})} \sup_{g \in G} \sum_{i=1}^N \sigma_i g(z_i)$$

Intuition 1: suppose I have drawn 2 data sets  $S_1, S_2$

$$\begin{aligned} \sup_g |\hat{R}_{S_1}(f) - \hat{R}_{S_2}(f)| &= \sup_{g \in G} \frac{1}{N} \left[ \sum_{(x,y) \in S_1} g(z_i) - \sum_{(x,y) \in S_2} g(z_i) \right] \\ &= \sup_g \frac{1}{N} \sum_{(x,y) \in S_1 \cup S_2} \sigma_i g(z_i) \end{aligned}$$

where  $\sigma_i = 1$  if  $(x,y) \in S_1$ ,  
 $-1$  otherwise.



Assume  $S$  is given,

we average  $\sup_g |\hat{R}_{S_1}(f) - \hat{R}_{S_2}(f)|$  over

all partitions of  $S$  in  $(S_1, S_2)$

we get Rademacher complexity

Intuition 2: Measure how well  $F$  can fit noisy labels.

Rademacher lemma:

$$\mathbb{E}_{S \sim P^N} [\text{UnRep}(F, S)] \leq 2 \mathbb{E}_{S \sim P^N} [\hat{Rad}(G)]$$

Then (PAC with Rademacher)

Assume  $|l(f(x), y)| \leq c$  for all  $(x, y)$

For all  $f \in \mathcal{F}$

if  $S \sim \mathcal{P}^N$ , with probab 1- $\delta$

$$R(f) - \hat{R}_S(f) \leq 2 \hat{R}_{\text{rad}_S}(\mathcal{F}) + 4e \sqrt{\frac{2 \ln 4\delta}{N}}$$

we conclude the PAC result:

$$R(f) - R(f_F) \leq ? \quad (\text{exercise})$$

exercise: 1) let  $\mathcal{G} = \{z \mapsto \alpha : \alpha \in [-1, 1]\}$   
what is  $\hat{R}_{\text{rad}_S}(\mathcal{G}) = ?$

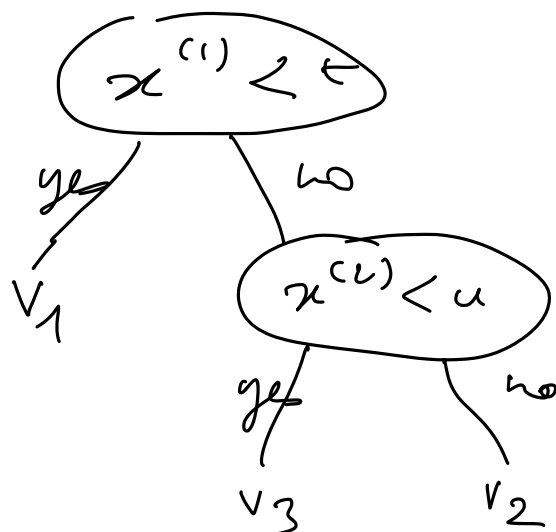
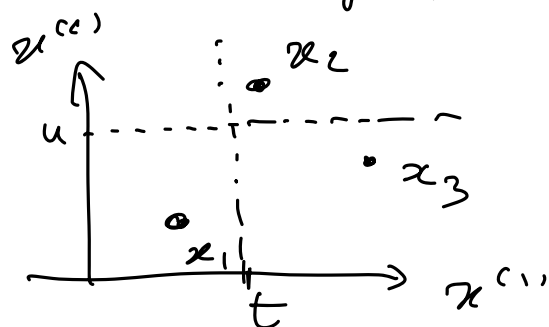
$$\hat{R}_{\text{rad}_S}(\mathcal{G}) = \frac{1}{N} \mathbb{E}_{\sigma_1 \dots \sigma_N \sim \text{unif}(-1, 1)} \sup_{g \in \mathcal{G}} \sum_{i=1}^N \sigma_i g(z_i)$$

2) let  $\mathcal{G} = \{ \text{decision trees which can output } \pm 1 \text{ at the leaves} \}$   
 $\hat{R}_{\text{rad}_S}(\mathcal{G}) = ?$

$$1) \sup_{\alpha \in [-1, 1]} \sum_{i=1}^N \sigma_i \alpha = \sup_{\alpha \in \{-1, 1\}} \sum_{i=1}^N \sigma_i \alpha = \left| \sum_{i=1}^N \sigma_i \right|$$

$$\begin{aligned} \hat{R}_{\text{rad}_S}(\mathcal{G}) &= \frac{1}{N} \mathbb{E}_{\sigma_1 \dots \sigma_N} \left| \sum_{i=1}^N \sigma_i \right| = \frac{1}{N} \mathbb{E}_{\sigma_1 \dots \sigma_N} \sqrt{(\sum \sigma_i)^2} \\ &\leq \frac{1}{N} \sqrt{\mathbb{E}_{\sigma_1 \dots \sigma_N} (\sum_{i=1}^N \sigma_i)^2} \\ &= \frac{1}{N} \sqrt{\text{Var}(\sum \sigma_i)} \quad (\text{Jensen inequality}) \\ &= \frac{1}{N} \sqrt{N \text{Var}(\sigma_1)} \\ &= \frac{1}{\sqrt{N}} \end{aligned}$$

2) let us choose a dataset of 3 points.



$\sigma_1$	$\sigma_2$	$\sigma_3$	$g$	$\sum \sigma_i g(x_i)$
-1	-1	-1	$T(-1, -1, -1)$	3
-1	-1	1	$T(-1, -1, 1)$	3
-1	1	-1	$T(-1, 1, -1)$	3
-1	1	1	$T(-1, 1, 1)$	$\vdots$
1	-1	-1	$T(1, -1, -1)$	$\vdots$
1	-1	1	$\vdots$	$\vdots$
1	1	-1	$\vdots$	$\vdots$
1	1	1	$\vdots$	$\vdots$

Call  $T(j, k, l)$   
the above tree where  
 $v_1 = j, v_2 = k, v_3 = l$

Here,  $\forall \sigma_1, \dots, \sigma_N, \sup_g \sum_{i=1}^N \sigma_i g(x_i) = N$   
 $\hat{\text{Rad}}_S(G) = 1.$

$\Rightarrow$  in the above then this bound is completely useless!

Thm: for any  $P$ , for  $F = \{x \mapsto x^\top \theta : \|\theta\|_2 \leq W_2\}$   
 for any 1-Lipschitz loss (hinge, logistic loss...)

$$\text{UnRep}(F, S) \leq \frac{W_2 X_2}{\sqrt{N}} + 4 X_2 \sqrt{\frac{2}{N} \ln \frac{2}{\delta}}$$

Rademacher  
complexity

$$X_2 = \sup_{x \in X} \|x\|_2$$