

# LECTURE 27

## EMPIRICAL PROCESSES

- Let  $X$  be a r.v. taking values in  $\mathcal{X}$  set  $\mathcal{X}$ ,  
 $X_1, \dots, X_n$  be independent copies of  $X$ .
- Law of large numbers  $\Rightarrow \forall$  Boolean function  $f: \mathcal{X} \rightarrow \{0,1\}$ :

$$\frac{1}{n} \sum_1^n f(x_i) \xrightarrow{\text{a.s.}} \mathbb{E}f(x) \quad \text{as } n \rightarrow \infty$$

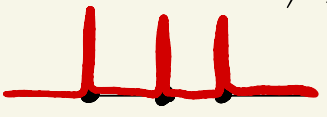
• Deviation:

$$\begin{aligned} \mathbb{E} \left| \frac{1}{n} \sum_1^n f(x_i) - \mathbb{E}f(x) \right| &\leq \left( \mathbb{E} \left| \dots \right|^2 \right)^{1/2} = \text{Var} \left( \frac{1}{n} \sum_1^n f(x_i) \right)^{1/2} \\ &= \left[ \frac{1}{n^2} \sum_1^n \underbrace{\text{Var}(f(x_i))}_{\substack{\wedge \\ \downarrow \\ \text{1 (f Boolean)}}} \right]^{1/2} = \frac{1}{\sqrt{n}} \quad \text{("Weak LLN")} \end{aligned}$$

• Is this true uniformly over all Boolean functions  $f$ ?

$$\mathbb{E} \sup_{\forall f} \left| \frac{1}{n} \sum_1^n f(x_i) - \mathbb{E}f(x) \right| \rightarrow 0 \quad \text{as } n \rightarrow \infty ?$$

**NO** For  $f = \mathbb{1}_{\{x_1, \dots, x_n\}}$ ,  $\frac{1}{n} \sum_1^n f(x_i) = 1$  but  $\mathbb{E}f(x) = 0$



• But it is true uniformly over  $f \in \mathcal{F}$  whenever  $vc(\mathcal{F}) < \infty$ :

### VLM (Uniform Law of Large Numbers)

If  $\mathcal{F}$  is a Boolean class with  $d = vc(\mathcal{F}) < \infty$ , then

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}f(x) \right| \leq C \sqrt{\frac{d \log n}{n}}$$

Remarks ① Same rate  $O(1/\sqrt{n})$  as in WLLN!

②  $\left( \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}f(x) \right)_{f \in \mathcal{F}}$  is called an *empirical process*.

The proof uses a new tool:

The proof of the Uniform Law of Large Numbers is based on:

Symmetrization Lemma [Gine-Zinn]

Let  $X_1, X_2, \dots, X_n$  be iid r.v.'s;  
 let  $\varepsilon_1, \dots, \varepsilon_n$  be independent  $\pm 1$  with prob  $1/2$ .  
 Then 
$$\mathbb{E} \left| \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \leq 2 \mathbb{E} \left| \sum_{i=1}^n \varepsilon_i X_i \right|$$

Proof

$$\mathbb{E} \left| \sum_i X_i - \mathbb{E} \left( \sum_i X_i \right) \right| = \mathbb{E}_X \left| \sum_i X_i - \mathbb{E}_{X'} \left( \sum_i X_i' \right) \right|$$

where  $X_i'$  are independent copies of  $X_i$

$$= \mathbb{E}_X \left| \mathbb{E}_{X'} \left( \sum_i (X_i - X_i') \right) \right| \leq \mathbb{E}_X \mathbb{E}_{X'} \left| \sum_i \underbrace{(X_i - X_i')}_{\text{dist}} \right|$$

Jensen ineq.

by symmetry:

$$X_i - X_i' \stackrel{\text{dist}}{=} X_i' - X_i$$

$$= \mathbb{E}_X \mathbb{E}_{X'} \mathbb{E}_\varepsilon \left| \sum_i \varepsilon_i (X_i - X_i') \right| \leq \mathbb{E} \left| \sum_i \varepsilon_i X_i \right| + \mathbb{E} \left| \sum_i \varepsilon_i X_i' \right|$$

!!  
E

$$\sum_i \varepsilon_i X_i + \sum_i \varepsilon_i X_i'$$

same distr.

$$= 2 \mathbb{E} \left| \sum_i \varepsilon_i X_i \right| \quad \text{QED.}$$

## Proof of ULLN

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right| = \frac{1}{n} \mathbb{E} \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n (f(X_i) - \mathbb{E} f(X_i)) \right|$$

$$\leq 2 \cdot \frac{1}{n} \mathbb{E} \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \varepsilon_i f(X_i) \right|$$

by a "uniform" version of Symmetrization Lemma (DIY)

"Rademacher complexity" of  $\mathcal{F}$

- Condition on r.v's  $X_i$  (treat them as fixed #'s). Randomness remains in  $\varepsilon_i$ .

Use Hoeffding inequality:

$$\forall f \in \mathcal{F},$$

$$P \left\{ \frac{1}{n} \left| \sum_{i=1}^n \varepsilon_i \underbrace{f(x_i)}_{a_i} \right| \geq t \right\} \leq 2 \exp \left( - \frac{(nt)^2}{2 \sum_{i=1}^n a_i^2} \right) \leq 2 \exp \left( - \frac{nt^2}{2} \right).$$

$a_i \in \{0, 1\}$  (Boolean class)

- Restriction:  $\sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right|$  is determined by the values of  $f \in \mathcal{F}$  on the sample  $x_1, \dots, x_n$  only.

$\Rightarrow$  we can replace the class  $\mathcal{F}$  by its restriction

$$\mathcal{F}_n := \left\{ f|_{\{x_1, \dots, x_n\}} : f \in \mathcal{F} \right\}.$$

- Sauer-Shelah Lemma  $\Rightarrow$

$$|\mathcal{F}_n| \leq \sum_{k=1}^d \binom{n}{k} \leq (en)^d \quad \text{where } d = \text{vc}(\mathcal{F}).$$

- Union Bound:

$$P \left\{ \sup_{f \in \mathcal{F}} \left| \frac{2}{n} \sum_{i=1}^n \varepsilon_i f(x_i) \right| \geq t \right\} \leq |\mathcal{F}_n| \cdot 2 \exp \left( - \frac{t^2 n}{2} \right)$$

$$\leq 2 \exp \left( d \log(en) - \frac{t^2 n}{2} \right)$$

$$\leq 2 \exp \left( - \frac{t^2 n}{4} \right)$$

$$\forall t \geq 10 \sqrt{\frac{d \log n}{n}} =: t_0$$

- Integral identity:

$$E Z = \int_0^\infty P\{Z \geq t\} dt \leq \int_0^{t_0} P\{Z \geq t\} dt + \int_{t_0}^\infty P\{Z \geq t\} dt$$

$$\leq t_0 + \int_{t_0}^\infty 2 \exp \left( - \frac{t^2 n}{4} \right) dt \leq 2 t_0.$$

QED