

Study Group  
on

Empirical Processes.

Contact: [pasupath@purdue.edu](mailto:pasupath@purdue.edu)

- Apology for handwritten notes!
- How did we get here?
- Background and preparation
- Frequency and Other Logistics
- Agenda

# Why Study Empirical Processes?

## Four motivating contexts

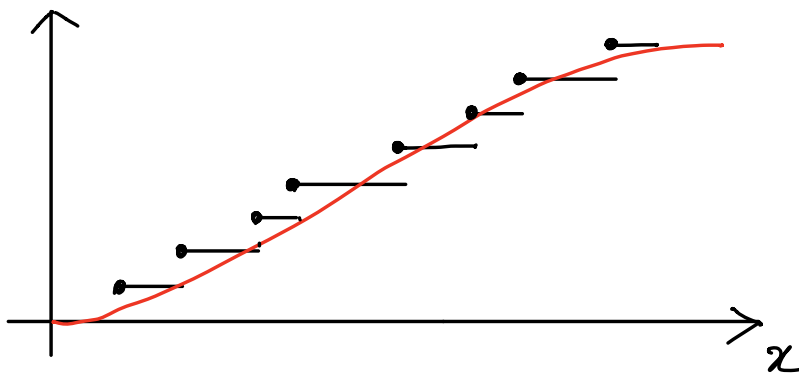
- I Estimate a cdf or a quantile function.
- II Stochastic Optimization & Stochastic Fixed ~~Point~~ Problem in Euclidean space.
- III Stochastic Optimization & Stochastic Fixed ~~Point~~ Problem in "non standard" spaces.
- IV Analyzing the bootstrap. X

## I. Estimate a CDF, Quantile

Let  $X_1, X_2, \dots$  be iid random variables with distribution function  $F$ .

The natural estimator of  $F$  is:

$$F_n(x) := \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{(-\infty, x]}(X_i)$$



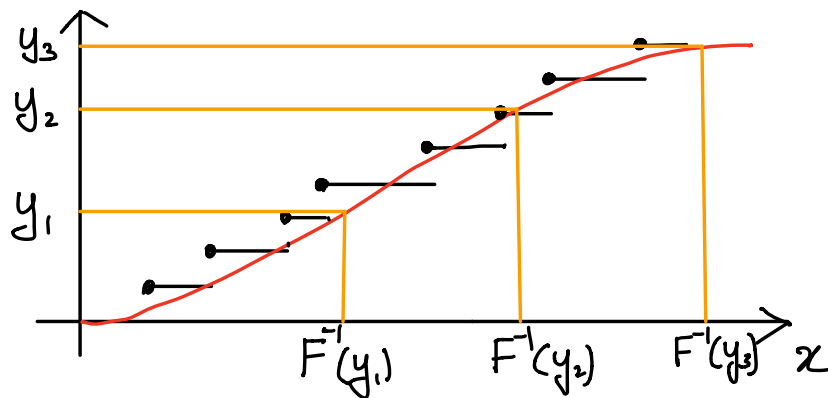
## I. Estimate a CDF, Quantile

The natural estimator of the quantile function

$$F^{-1}(y) := \inf\{x: F(x) \geq y\}$$

is the sample quantile function

$$Q_n(y) := \inf\{x: F_n(x) \geq y\}.$$



## I. Estimate a CDF, Quantile

What can we say about  $F_n$   
as an estimator of  $F$ ?

Specifically:

- (i)  $\sup_x |F_n(x) - F(x)| \xrightarrow{a.s.} 0$ ?
- (ii) How fast?

Analogous questions about  $Q_n$   
as an estimator of  $F^{-1}$ .

## I. Estimate a CDF, Quantile

The principal route to answering such questions

analyzes the empirical process:

$$n^{\frac{1}{2}} (P_n(x) - P(x))$$

$$\beta_n(x) = \sqrt{n} (F_n(x) - F(x)), -\infty < x < \infty$$

$$q_n(y) = \sqrt{n} (Q_n(y) - F^{-1}(y)), 0 < y < 1$$

## I. Estimate a CDF, Quantile

For example, if  $X_j \in [0, 1]$ ,  $j=1, 2, \dots, n$   
then,

$$\beta_n \Rightarrow B \quad (1)$$

where  $B$  is the Gaussian random element of  $\mathcal{D}[0, 1]$  specified by  
 $E[B(t)] = 0$ ,  $E[B(s)B(t)] = F(s \wedge t) - F(s)F(t)$ .  
(Theorem 14.3 in Billingsley, 1999.)

Beware: (1) is weak convergence in metric space.



## I. Estimate a CDF, Quantile

So what?

### Weak Invariance Principle

Because of (1), it may be expected that "operations" on  $\beta_n$  will "behave like operations" of  $B$ .

Hence the behavior of  $F_n$  "can be approximated" by passing to the limit.

## I. Estimate a CDF, Quantile

Two Classic Examples (DasGupta, 2011)

(i) Kolmogorov-Smirnov Statistic

$$\sup_{0 < x < 1} |B_n(x)| \implies \sup_{0 < x < 1} |B(x)|$$

where  $\sup_{0 < x < 1} |B(x)|$  has the

known cdf  $1 - 2 \sum_{k=1}^{\infty} (-1)^{k-1} e^{-2k^2 y^2}, y > 0$

(ii) Cramér-von Mises

$$\int_0^1 B_n^2(x) dF \implies \int_0^1 B^2(x) dx$$

## I. Estimate a CDF, Quantile

The analogous setting for the quantile process

Suppose  $F$  is absolutely continuous with density  $f > 0$ ,  $f$  differentiable and  $F(x)(1-F(x)) \frac{|f'(x)|}{f^2(x)}$  is uniformly bounded on  $\text{supp.}(f)$ . Then, for any  $c \in (0,1)$

$$\sup_{c \leq y \leq 1-c} |q_n(y) f(F^{-1}(y))| \implies \sup_{c \leq y \leq 1-c} |B(y)|$$

## I. Estimate a CDF, Quantile

Stepping back...

$$B_n(x) = \sqrt{n} \left( \underbrace{F_n(x) - F(x)}_{\text{"est. error"}}, x \in \mathbb{R} \right)$$

"scaling"                      "parameter"

$$\Rightarrow \left\{ B(x), x \in \mathbb{R} \right\} \rightarrow \text{Gaussian process}$$

**Notice:** We can also write the above as

$$B_n(x) = \sqrt{n} \left( P_n((-\infty, x]) - P((-\infty, x]) \right), x \in \mathbb{R}$$
$$\Rightarrow \left\{ P_B((-\infty, x]), x \in \mathbb{R} \right\}$$

## II Stochastic Optimization

### Example 1

Relate employment with education.

Data  $(Y_i, Z_i)$ ,  $i=1, 2, \dots, n$

where

$$Y_i = \begin{cases} 1 & \text{if } i \text{ has a job.} \\ 0 & \text{otherwise} \end{cases}$$

$Z_i$  is the education of individual  $i$ .

Assume:

$$P(Y=1 \mid Z=z) = G(\theta^* z), \theta^* \in \mathbb{R}^d$$

$$\text{where } G(\xi) = (1 + e^{-\xi})^{-1}, \xi \in \mathbb{R}.$$

## II Stochastic Optimization

Optimization Problem (Classic MLE)

Find  $\theta_n \in \operatorname{argmax} H_n(\theta)$  — (P<sub>n</sub>)

where

$$H_n(\theta) := \frac{1}{n} \sum_{j=1}^n H(\theta, Y_j)$$

$$H(\theta, Y) := \log G_1(\theta z)^Y (1 - G_1(\theta z))^{1-Y}$$

$$\theta^* := \operatorname{argmax} h(\theta) = \mathbb{E}[H(\theta, Y)]$$

## II Stochastic Optimization

### Key Questions

- (i) Does  $\theta_n$  converge to  $\theta$  in any sense, e.g.,  $\|\theta_n - \theta^*\| \xrightarrow{\text{a.s.}} 0$ ?  
 $H(\theta_n) - H(\theta^*) \xrightarrow{\text{a.s.}} 0$ ?
- (ii) How fast?

## II Stochastic Optimization

In answering questions in A. & B.

EP is the natural viewpoint. Why?

Define the EP

$$\varepsilon_n(\theta) = \sqrt{n} \left( H_n(\theta) - h(\theta) \right), \theta \in \mathbb{R}^d$$

parameter

scaling

error

Principle: If  $\varepsilon_n$  can be shown to converge to a Gaussian process, then, answers to (i) and (ii) become possible.



## II Stochastic Optimization

There is nothing special about the MLE example.

General Stochastic Optimization:

Find  $\operatorname{argmin}_{\theta \in \Theta} h(\theta) := \mathbb{E}[H(\theta, Y)]$

$$\theta \in \Theta \subseteq \mathbb{R}^d$$

where  $Y$  in  $(\mathcal{Y}, \mathcal{A})$ ,  $H: \mathbb{R}^d \times \mathcal{Y} \rightarrow \mathbb{R}$ .

(50)

## II Stochastic Optimization

Find  $\operatorname{argmin}_{\theta \in \Theta \subseteq \mathbb{R}^d} h(\theta) := \mathbb{E}[H(\theta, Y)]$

$$\theta \in \Theta \subseteq \mathbb{R}^d$$

where  $Y$  in  $(\mathcal{Y}, \mathcal{A})$ ,  $H: \mathbb{R}^d \times \mathcal{Y} \rightarrow \mathbb{R}$ .

"Sample Version" of (SO).

Find  $\operatorname{argmin}_{\theta \in \Theta \subseteq \mathbb{R}^d} H_n(\theta) := \frac{1}{n} \sum_{j=1}^n H(\theta, Y_j)$

$$\theta \in \Theta \subseteq \mathbb{R}^d$$

where  $Y_j, j=1, 2, \dots \in (\mathcal{Y}, \mathcal{A})$ , iid.

## II Stochastic Optimization

Example Settings for SO.

- Virtually all of regression
- A large fraction of optimization problems in Operations Research, Computer Science, & Engineering.

### III Stochastic Optimization

The non-Euclidean context emerges...

Example 2. Recall that in Example 1,

$$P(Y=1 | Z=z) = G_1(\theta^* z), \theta^* \in \mathbb{R}^d$$

$$\operatorname{argmax}_{\theta \in \mathbb{R}^d} h(\theta) := \mathbb{E}[H(\theta, Y)]$$

What if, instead:

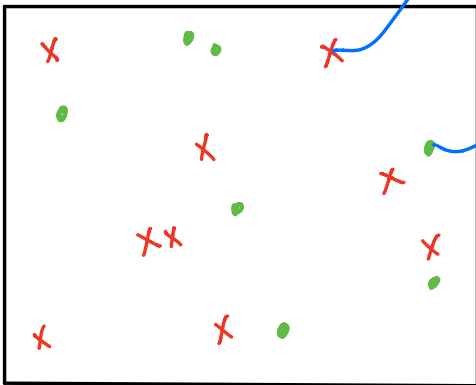
$$P(Y=1 | Z=z) = G_0(z)$$

$$\operatorname{argmax}_{G_1 \in \mathcal{G}} h(G_1) := \mathbb{E}[H(G_1, Y)]$$
$$G_1 \in \mathcal{G} = \{\mathbb{R} \rightarrow [0, 1], \text{increasing}\}$$

### III Stochastic Optimization

#### Example 3

S



→ "event" occurs according to prob. measure  $\Lambda$ .

→ "good guys" sprinkled according to Poisson measure  $\Theta$

exp. prob. of death

$$\begin{aligned}
 h(\theta) &= \int_S \Lambda(dx) \int_0^1 P(f(R_x, \theta) \geq u) du \\
 &= \int_S \Lambda(dx) \int_0^1 \exp\{-\theta(B(x, \sqrt{f^{-1}(u)}))\} du
 \end{aligned}$$

Find  $\operatorname{argmin} h(\theta)$ ,  $\int \theta(dx) = b$

### III Stochastic Optimization

Define the EP

$$\varepsilon_n(\theta) = \sqrt{n} \left( \underbrace{H_n(\theta)}_{\text{parameter}} - \underbrace{h(\theta)}_{\text{error}} \right), \theta \in \Theta$$

scaling

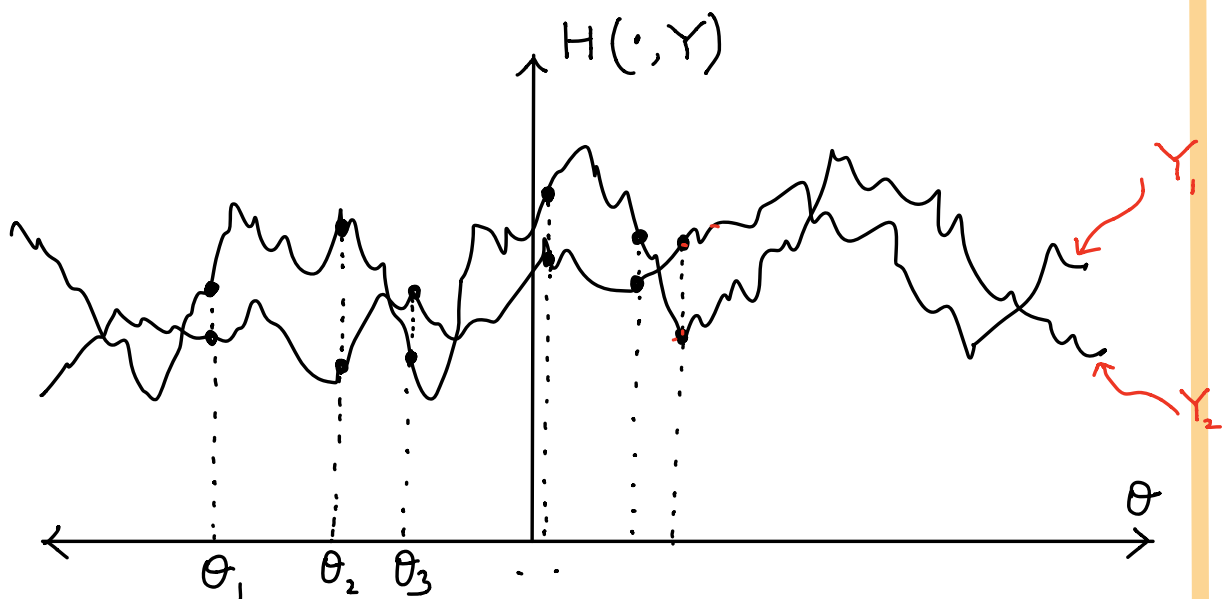
Q.1 What are conditions for  
 $\sup_{\theta \in \Theta} |H_n(\theta) - h(\theta)| \xrightarrow{\text{a.s.}} 0$ ?

Q.2. At what rate?

### III Stochastic Optimization

Answers to Q.1 and Q.2 will be stated in terms of **entropy**.

or "diversity" or "richness" of the class  $\{H(\theta, Y_j), \theta \in \Theta, j=1, 2, \dots, n\}$ .



## III Stochastic Optimization

### AGENDA

#### I. Basic Results

- (i) Donsker's thm for partial sums
- (ii) EP for CDFs
- (iii) Quantile process

#### II. M-Estimation

- (i) Entropy
- (ii) Uniform convergence
- (iii) Asymptotic equicontinuity.



## WEAK CONVERGENCE ESSENTIALS

Today's Agenda:

Demonstrate the first example of weak convergence in "non-Euclidean" space.

- Weak convergence essentials
- Specialize to  $C[0,1]$ , main theorem
- Wiener measure, Donsker's Thm.

## WEAK CONVERGENCE ESSENTIALS

Our treatment is on a metric space  $S = (S, \rho)$ , where  $S$  is a set and  $\rho$  is a metric on  $S$ .

For  $x, y, z \in S$

$$(M1) \quad \rho(x, y) \in [0, \infty)$$

$$(M2) \quad \rho(x, y) = 0 \text{ iff } x = y$$

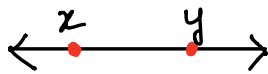
$$(M3) \quad \rho(x, y) = \rho(y, x)$$

$$(M4) \quad \rho(x, z) \leq \rho(x, y) + \rho(y, z)$$

(See D2011 or K1978)

# WEAK CONVERGENCE ESSENTIALS

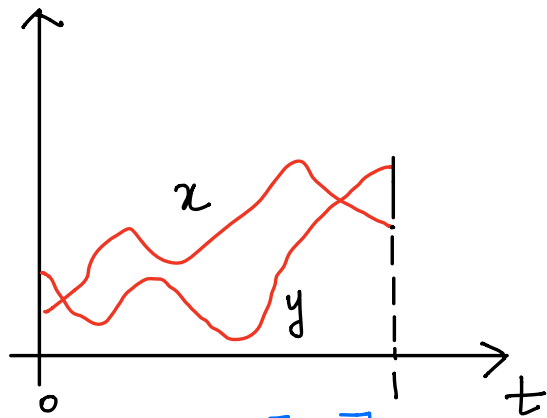
## I Reals ( $\mathbb{R}$ )



$$S = \mathbb{R}$$

$$f(x, y) = |x - y|$$

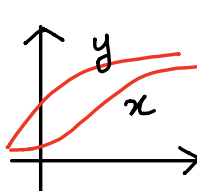
## II Function Space



$$S = C[0, 1]$$

$$f(x, y) = \sup_{0 \leq t \leq 1} |x(t) - y(t)|$$

## IV Measure Space ( $M$ )



$x, y \in M$  with  
cdf  $F_x, F_y$

$$f(x, y) = \sup_{z \in \mathbb{R}^d} |F_x(z) - F_y(z)|$$

(See D2011)

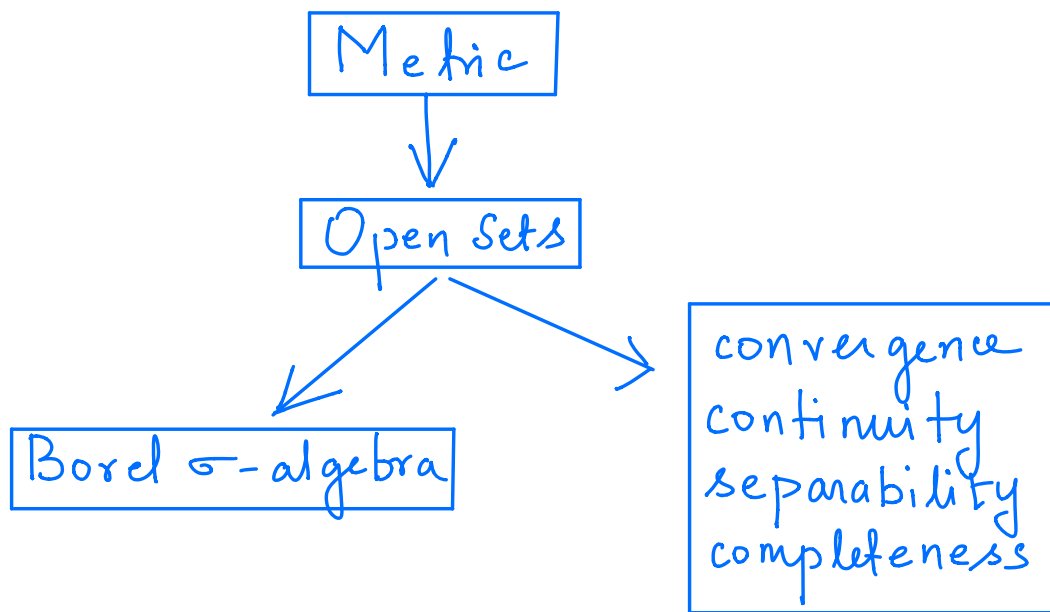
## III Sequence Space, $S = l^\infty$ .

$$x = (\xi_1, \xi_2, \dots); y = (\eta_1, \eta_2, \dots)$$

$$f(x, y) = \sup_{j \in \mathbb{N}} |\xi_j - \eta_j|$$

# WEAK CONVERGENCE ESSENTIALS

The metric "does a lot of work."



(Sometimes, the Borel  $\sigma$ -algebra is "too big.")

So, defining the metric is key to defining the measurable space  $(S, \mathcal{S})$ .

## WEAK CONVERGENCE ESSENTIALS

A probability measure on  $\mathcal{S}$  is a non-negative, countably additive set function with  $P(S) = 1$ .

$P_n$  converges weakly to  $P$  means, for  $A \in \mathcal{S}$

$$P_n(A) \rightarrow P(A) \text{ if } P(\partial A) = 0.$$

Notation:

$$P_n \Rightarrow P$$

## WEAK CONVERGENCE ESSENTIALS

What is the relevance of  
the  $P$ -continuity set condition  
 $P(\partial A) = 0$ ?

### Understand by Analogy

In  $\mathbb{R}$ ,  $X_n \Rightarrow X$  means

$F_{X_n}(t) \rightarrow F_X(t)$  at points  $t$

where  $F_X$  is continuous.

$$F_{X_n}(t) = P_n((-\infty, t]); \quad F_X(t) = P((-\infty, t])$$

## WEAK CONVERGENCE ESSENTIALS

Let  $X: \Omega \rightarrow S$  be a mapping from  $(\Omega, \mathcal{F}, P)$  to the metric space  $S$ .

$X$  is measurable  $\mathcal{F}/\mathcal{B}$ . ( $X \in \mathcal{F}/\mathcal{B}$ )

Distribution of  $X$  is the probability measure induced by  $X$ .

$P = P X^{-1}$ , that is:

$$P(A) = (P X^{-1})(A) = P(X^{-1}(A))$$

$$= P(\omega: X(\omega) \in A)$$

## WEAK CONVERGENCE ESSENTIALS

(i) - (iv) mean the same.

$$(i) \quad X_n \Rightarrow X$$

$$(ii) \quad P_n \Rightarrow P$$

$$(iii) \quad X_n \Rightarrow P$$

$$(iv) \quad P_n \Rightarrow X$$

And, equivalent to:

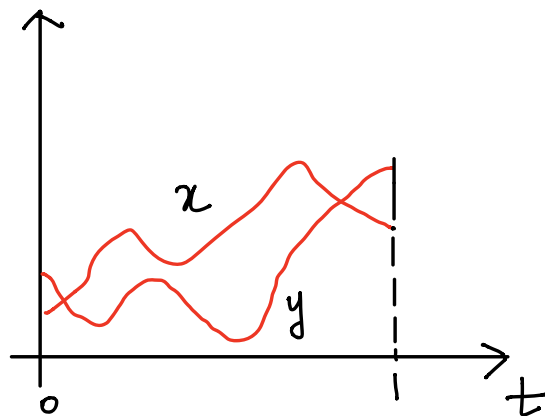
$$E[f(X_n)] \rightarrow E[f(X)]$$

for all bounded, uniformly continuous  $f$ .



## WEAK CONVERGENCE ESSENTIALS

Let's live in  $C[0,1]$  today.



$C = C[0,1]$  is the space of continuous functions endowed with the "uniform topology":

$$p(x, y) = \|x - y\| = \sup_t |x(t) - y(t)|$$

$(C, \mathcal{C})$ : corresponding measurable space.

## WEAK CONVERGENCE ESSENTIALS

Random Function; Random Vector; Random Variable.

$$(\Omega, \mathcal{F}, \mathbb{P}) \text{ and } X: \Omega \rightarrow \mathbb{C}[0, 1]$$

If  $X \in \mathcal{F}/\mathbb{C}$ , then it is called a  
random function.

random variable

$$\pi_t X(\omega) := X(t, \omega), \quad \pi_t X \in \mathcal{F}/\mathbb{R}$$

$$\pi_{t_1, t_2, \dots, t_k} X(\omega) := (X(t_1, \omega), X(t_2, \omega), \dots, X(t_k, \omega))$$

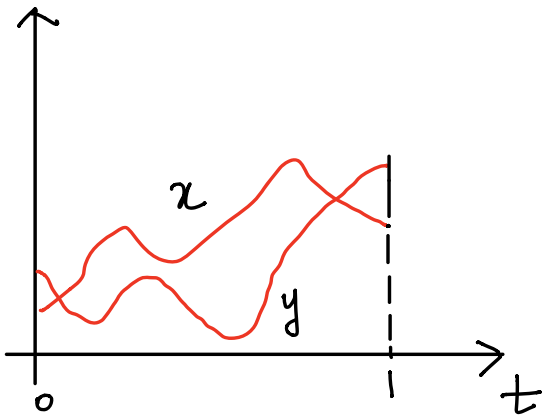
$$\pi_{t_1, t_2, \dots, t_k} X \in \mathcal{F}/\mathbb{R}^k$$

random vector

# WEAK CONVERGENCE ESSENTIALS

## Modulus of Continuity.

$$m_x(\delta) := \sup_{|s-t| \leq \delta} |x(s) - x(t)| \quad 0 < \delta \leq 1$$



" $m_x(\cdot)$  quantifies how  $\varepsilon$  changes with  $\delta$  in the  $\varepsilon$ - $\delta$  definition of continuity."

## Examples

(i) If  $x$  is  $L$ -Lipschitz,  $m_x(\delta) \leq L\delta$

(ii) If  $x$  is Hölder,  $m_x(\delta) \leq L\delta^\alpha$

## WEAK CONVERGENCE ESSENTIALS

Suppose  $X, X^1, X^2, \dots$  are random functions.

Theorem 7.5

iff

$$(X_{t_1}^n, X_{t_2}^n, \dots, X_{t_k}^n) \Rightarrow (X_{t_1}, X_{t_2}, \dots, X_{t_k})$$

$$\forall t_1, t_2, \dots, t_k \in [0, 1]$$

(FD)

and

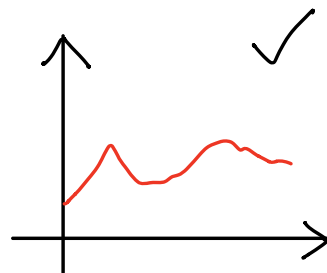
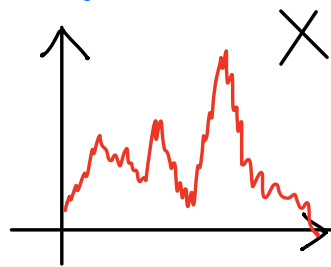
$$\lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} P(m(X^n, \delta) > \varepsilon) = 0,$$

(KC)

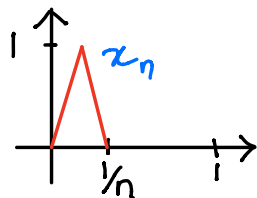
then  $X^n \Rightarrow X$ .

## WEAK CONVERGENCE ESSENTIALS

(Kc) means random function changes can be "controlled."



Ex. 1 (FD) is not enough:



$x_n \rightarrow x = 0$  and hence  $\delta_{x_n} \not\rightarrow \delta_x$ .

Ex. 2 (FD) is enough:  $S = \mathbb{R}^\infty$

## WEAK CONVERGENCE ESSENTIALS

(Proof Sketch of 7.5)

- (FD) implies  $P_n \pi_0^{-1}$  is tight.
- $\{P_n \pi_0^{-1}\}$  tight + KC  $\Leftrightarrow \{P_n\}$  tight.
- FD +  $\{P_n\}$  tight is sufficient.

## WEAK CONVERGENCE ESSENTIALS

(KC) looks "artificial."

Is it too much?

In short, because we are in  $C[0, 1]$ , the answer is "No."

(KC) is fundamental and controls "richness" to just right extent.

## WIENER MEASURE

Wiener measure,  $W$ , is a probability measure on  $(C, \mathcal{C})$  having two properties.

$$(A) \quad W[x_t \leq \alpha] = \frac{1}{\sqrt{2\pi t}} \int_{-\infty}^{\alpha} \exp\left\{-\frac{1}{2} \frac{u^2}{t}\right\} du$$

(B) for  $0 \leq t_0 \leq t_1 \leq \dots \leq t_k = 1$ ,  
 $x_{t_1} - x_{t_0}, x_{t_2} - x_{t_1}, \dots, x_{t_k} - x_{t_{k-1}}$   
are independent under  $W$ .



# WIENER MEASURE

## Two Crucial Things.

1.  $W \left[ x_{t_i} - x_{t_{i-1}} \leq \alpha_i, i=1, 2, \dots, k \right]$

$$= \prod_{i=1}^k \frac{1}{\sqrt{2\pi(t_i - t_{i-1})}} \int_{-\infty}^{\alpha_i} \exp\left\{-\frac{1}{2} \frac{u^2}{t_i - t_{i-1}}\right\}$$

Thus, the finite-dimensional distributions are specified.

$$\left( \int_{x_{t_1}, x_{t_2}, \dots, x_{t_k}} f(z_1, z_2, \dots, z_k) = \int_{z_{t_1}} f(z_1) \int_{z_{t_2} - z_{t_1}} f(z_2 - z_1) \dots \int_{z_{t_k} - z_{t_{k-1}}} f(z_k - z_{k-1}) \right)$$

2. The existence of  $W$  needs to be proved: at most one, sometimes none.

## DONSKER'S THEOREM

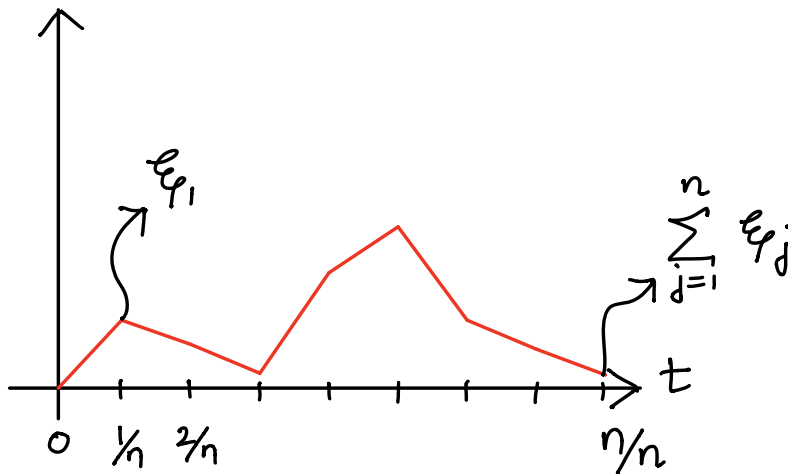
Lets assume  $W$  exists, and construct a sequence  $\{X^n\}$  such that

$$X^n \Rightarrow W.$$

Let  $\xi_1, \xi_2, \dots$  be i.i.d random variables such that  $E[\xi_1] = 0$  and  $\text{Var}(\xi_1) = \sigma^2 \in (0, \infty)$ .

$$S_n = \begin{cases} 0 & n=0 \\ \xi_1 + \xi_2 + \dots + \xi_n, & n \geq 1 \end{cases}$$

# DONSKER'S THEOREM



$$X_t^n(\omega) = \frac{1}{\sigma\sqrt{n}} \left( S_{\lfloor nt \rfloor}(\omega) + (nt - \lfloor nt \rfloor) \xi_{\lfloor nt \rfloor + 1}(\omega) \right)$$

$t \in [0, 1]$

Notice that  $X^n \in C[0, 1]$

## DONSKER'S THEOREM

Due to Theorem 7.5, we would have proved

$$X^n \Rightarrow W$$

if we can show (FD) and (KC) hold.

# DONSKER'S THEOREM

Fix  $t$ :

$$X_t^n(\omega) = \underbrace{\frac{1}{\sigma\sqrt{n}} S_{\lfloor nt \rfloor}(\omega)}_{\substack{\Downarrow \\ \sqrt{t} N \\ \text{(Lindeberg-Levy CLT)}}} + \underbrace{\frac{(nt - \lfloor nt \rfloor) \xi_{\lfloor nt \rfloor + 1}(\omega)}{\sigma\sqrt{n}}}_{\substack{\downarrow \\ o_p(1) \\ \text{(by Chebyshev)}}$$

Fix  $s, t$  with  $s \leq t$ :



$$(X_s^n(\omega), X_t^n(\omega) - X_s^n(\omega))$$

$$= \frac{1}{\sigma\sqrt{n}} (S_{\lfloor ns \rfloor}(\omega), S_{\lfloor nt \rfloor}(\omega) - S_{\lfloor ns \rfloor}(\omega)) + o_p(1)$$

$$\Rightarrow (\sqrt{s} N_1, \sqrt{t-s} N_2)$$

## DONSKER'S THEOREM

Hence,

$$\begin{aligned} (X_s^n(\omega), X_t^n(\omega)) &= (X_s^n(\omega), X_s^n(\omega) + X_t^n(\omega) - X_s^n(\omega)) \\ &\Rightarrow (\sqrt{s} N_1, \sqrt{s} N_1 + \sqrt{t-s} N_2) \\ &\quad \text{(by mapping thm.)} \end{aligned}$$

Fix  $t_1, t_2, \dots, t_k$



$$\begin{aligned} (X_{t_1}^n(\omega), X_{t_2}^n(\omega) - X_{t_1}^n(\omega), \dots, X_{t_k}^n(\omega) - X_{t_{k-1}}^n(\omega)) \\ \Rightarrow (\sqrt{t_1} N_1, \sqrt{t_2 - t_1} N_2, \dots, \sqrt{t_k - t_{k-1}} N_k) \end{aligned}$$

This proves (FD).

## DONSKER'S THEOREM

We will not prove (KC)... sorry but see pp. 88-90 in B99.

### Theorem 8.2

If  $\xi_1, \xi_2, \dots$  are i.i.d with mean zero and variance  $\sigma^2 \in (0, \infty)$ , then

$$X^n \Rightarrow W.$$

## DONSKER'S THEOREM

What have we skipped?

- ① Existence of  $W$
- ② Proof that  $X^n, n \geq 1$  satisfies (KC).

See pp. 88-90 in B99 for both.



## WEAK CONVERGENCE ESSENTIALS

Is (KC) in Theorem 7.5 too much?

(Prohorov's Thm., pp. 58)

FD + RC  $\iff$  Weak Convergence

And,

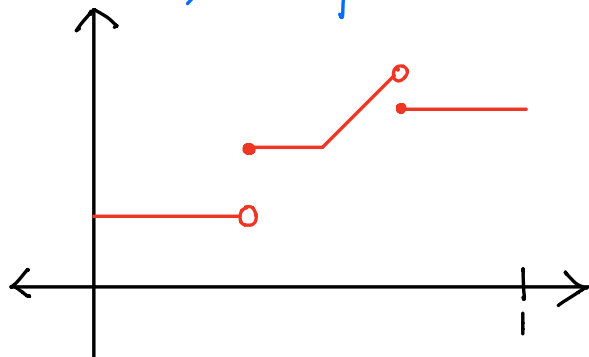
Under separability and completeness

FD + RC  $\iff$  FD + KC.

## WEAK CONVERGENCE ESSENTIALS

Lets go to  $\mathcal{D}[0, 1]$ .

Recall that  $\mathcal{D}[0, 1]$  is the class of cadlag functions, that is, right continuous functions with left limits, defined on  $[0, 1]$ .



# EMPIRICAL PROCESSES

## AGENDA

- (I) Step back ...
- (II) Generalization setup.
- (III) Entropy & Examples.
- (IV) ULLN (?)

①

Suppose  $X_1, X_2, \dots$  are i.i.d  
real-valued random variables  
with cdf  $F$ .

empirical cdf.

$$F_n(x) := n^{-1} \sum_{j=1}^n \mathbb{I}_{(-\infty, x]}(X_j), \quad x \in \mathbb{R}$$

$$Z_n(x) := \sqrt{n} (F_n(x) - F(x)), \quad x \in \mathbb{R}$$

empirical process (indexed by  $x$ )

(I)

(Glivenko-Cantelli, 1933)

$$\|F_n - F\|_\infty := \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \xrightarrow{\text{a.s.}} 0$$

(Donsker, 1952) std. Brownian Bridge  
on  $[0, 1]$ .

$$Z_n \Rightarrow Z \equiv \left\{ B(x), x \in [0, 1] \right\}$$

in  $\mathcal{D}(\mathbb{R}, \|\cdot\|_\infty)$ .

(I)

What are the analogues  
when  $X_1, X_2, \dots$  lie in a more  
general space  $\mathcal{X}$ , e.g.,  $C[0, 1]$ ,  
or Riemannian manifold?

How to define  $F(x) = P(X \leq x)$   
when  $X, x \in \mathcal{X}$ ?

Natural Idea :

(I)

$$F_n(x) = P_n((-\infty, x])$$

$$F(x) = P((-\infty, x])$$

The above suggests considering  $P_n(A)$  and  $P(A)$  for an appropriate class of sets  $\mathcal{C}$ ,  $\mathcal{C} \ni A \subseteq \mathcal{X}$ .

Now ask  $\sup_{A \in \mathcal{C}} |P_n(A) - P(A)| \xrightarrow{a.s.} 0?$

(I)

More conveniently, we can ask

$$\sup_{f \in \mathcal{F}} |P_n f - P f| \xrightarrow{\text{a.s.}} 0?$$

where  $\mathcal{F}$  is a class of real-valued functions having domain  $\mathcal{X}$ .

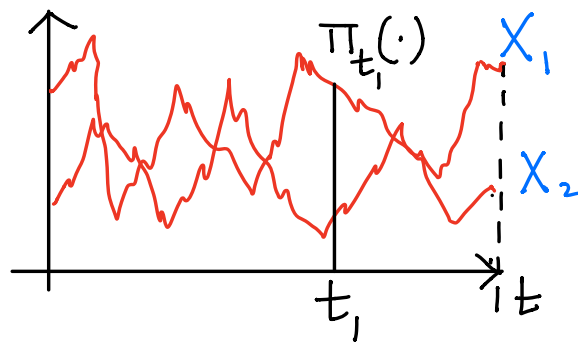
$$\left( P_n f = n^{-1} \sum_{j=1}^n f(X_j); P f = \int f dP \right).$$



# Two Examples.

(I)

(I)



$$\mathcal{X} \equiv C[0, 1]$$

$\mathcal{F} :=$  "projection functions"

$$\left\{ \pi_t(x) := x(t) : t \in [0, 1] \right\}$$

## Two Comments

Ⓡ

(II) Suppose  $\mathcal{X} \equiv \mathbb{R}$ .

Then

$\mathcal{F} :=$  "indicator functions"

$$= \left\{ \mathbb{I}_{(-\infty, x]}(\cdot), x \in \mathbb{R} \right\}$$

(I)

The "richness" or "complexity" of  $\mathcal{F}$  will determine the existence and nature of the GC and Donsker analogues.

II

Lets Setup...

II

$X_1, X_2, \dots \sim P$  are i.i.d in  $(\mathcal{X}, \mathcal{A})$

e.g.:  $\mathcal{X} = \mathbb{R}^d, C[0,1], \text{ etc.}$

### The Empirical Measure

$$P_n := n^{-1} \sum_{j=1}^n \delta_{X_j}$$

$$P_n(A) = n^{-1} \sum_{j=1}^n \mathbb{1}_A(X_j)$$

II

$P_n f, P_f$

For  $f : X \rightarrow \mathbb{R}$

$$P_f := \int f dP$$

$$P_n f := \int f dP_n = n^{-1} \sum_{j=1}^n f(X_j)$$

II

## Empirical Process

Suppose  $\mathcal{F}$  is a class of real-valued functions defined on  $\mathcal{X}$ .

$$\beta_n(f) := \sqrt{n} (P_n f - P f)$$

The stochastic process  $\{\beta_n(f), f \in \mathcal{F}\}$  is called an empirical process.

II

We hope to identify sufficient conditions on  $\mathcal{F}$  so that:

(Glivenko-Cantelli Analogue)

$$\|P_n - P\|_\infty^* = \left( \sup_{f \in \mathcal{F}} |P_n f - P f| \right)^* \xrightarrow{\text{a.s.}} 0$$

$\left\{ z: \mathcal{F} \rightarrow \mathbb{R}, \|z\|_\infty < \infty \right\}$

(Donsker Analogue)

$$\sqrt{n}(P_n - P) \Rightarrow G_1 \text{ in } \ell^\infty(\mathcal{F})$$

Where  $G_1$  is a  $P$ -Brownian bridge.



II

The sufficient conditions  
will be phrased in terms  
of some notion of the  
"complexity" of  $\mathcal{F}$ .

So, we now set up toward  
entropy of  $\mathcal{F}$ .

II

## $L_r(Q)$ norm

$Q$  is a measure on  $(X, \mathcal{A})$ .

For  $1 \leq r < \infty$ ,

$$\|f\|_{r,Q} := \left( \int |f|^r dQ \right)^{1/r}.$$

$\|f_1 - f_2\|_{r,Q}$  :  $L_r(Q)$  distance  
between  $f_1, f_2$ .

II

$\varepsilon$ -cover

$\{f_1, f_2, \dots, f_n\}$  is said to be

an  $\varepsilon$ -cover for  $\mathcal{F}$  if for any

$f \in \mathcal{F}$ ,  $\exists f_j$  such that

$$\|f - f_j\|_{r, Q} \leq \varepsilon.$$

( $f_1, f_2, \dots, f_n$  need not live in  $\mathcal{F}$ )

II

$\varepsilon$ -covering number.

$$N_n(\mathcal{F}, Q, \varepsilon) := \inf \left\{ n: \begin{array}{l} \exists \text{ an } L_n(Q) \text{ } \varepsilon\text{-cover} \\ \{f_1, f_2, \dots, f_n\} \text{ of } \mathcal{F} \end{array} \right\}$$

$\varepsilon$ -entropy (or metric entropy)

$$H_n(\mathcal{F}, Q, \varepsilon) := \log N_n(\mathcal{F}, Q, \varepsilon)$$

II

$\varepsilon$ -cover with bracketing.

$\left\{ [f_j^L, f_j^U]_{j=1}^n \right\}$  is said to be  
an  $\varepsilon$ -cover with bracketing for  $\mathcal{F}$   
if for each  $f \in \mathcal{F}$ :

$$\|f_j^U - f_j^L\|_{r, Q} \leq \varepsilon \quad \forall j.$$

$$\exists j \text{ s.t. } f_j^L \leq f \leq f_j^U$$

II

## $\varepsilon$ -covering with bracketing

$$N_{r,B}(\mathcal{F}, Q, \varepsilon) := \inf \left\{ n : \begin{array}{l} \exists \text{ an } \varepsilon\text{-cover with brae.} \\ \left\{ [f_j^L, f_j^U]_{j=1}^n \right\} \end{array} \right\}$$

## $\varepsilon$ -entropy with bracketing

$$H_{r,B}(\mathcal{F}, Q, \varepsilon) := \log N_{r,B}(\mathcal{F}, Q, \varepsilon)$$

II

## $\varepsilon$ -entropy for the sup. norm

$$\|f\|_{\infty} := \sup_{x \in X} |f(x)|$$

$$H_{\infty}(\mathcal{F}, \varepsilon) := \log N_{\infty}(\mathcal{F}, \varepsilon)$$

where

$$N_{\infty}(\mathcal{F}, \varepsilon) :=$$

$$\inf \left\{ n : \begin{array}{l} \exists \{f_1, f_2, \dots, f_n\} \text{ s.t.} \\ \sup_{f \in \mathcal{F}} \min_{1 \leq j \leq n} \|f - f_j\|_{\infty} \leq \varepsilon \end{array} \right\}$$

Notice: no dependence on  $Q$

## Comments

II

-  $Q$  need not be a prob. measure

-  $\|f\|_{p,Q} \uparrow p$  but

$$\lim_{p \rightarrow \infty} \|f\|_{p,Q} \neq \|f\|_{\infty} = \sup_{x \in X} |f(x)|$$

$\underbrace{\hspace{10em}}_{Q\text{-indep.}}$

-  $H_p(\mathcal{F}, Q, \varepsilon) \leq H_{p,B}(\mathcal{F}, Q, \varepsilon) \quad \forall \varepsilon > 0$

$H_{p,B}(\mathcal{F}, Q, \varepsilon) \leq H_{\infty}(\mathcal{F}, \frac{\varepsilon}{2})$  if

$Q$  is a prob. measure.



Entropy calculation examples. (III)

Example 1 (finite support)

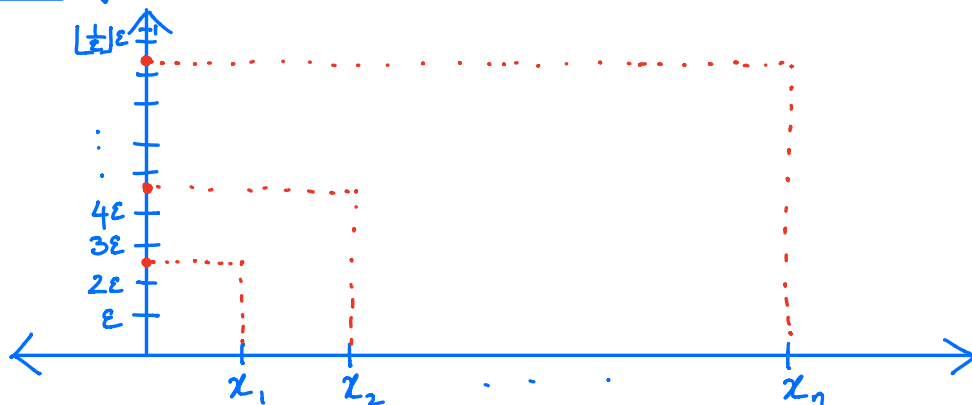
$$\mathcal{F} := \left\{ \begin{array}{l} \text{increasing functions} \\ f: X \subset \mathbb{R} \rightarrow [0, 1] \text{ where} \\ |X| = n < \infty. \end{array} \right\}$$

Then:

$$H_{\infty}(\mathcal{F}, \varepsilon) \leq \frac{1}{\varepsilon} \log \left( n + \frac{1}{\varepsilon} \right)$$

$$\forall \varepsilon > 0.$$

# Proof Sketch



Suppose  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$

$\mathcal{F} := \left\{ \begin{array}{l} \text{non-decreasing functions } f \\ \text{on } \mathcal{X} \text{ such that:} \\ f(x_j) \in \{i\epsilon, 0 \leq i \leq \lfloor L/\epsilon \rfloor\} \end{array} \right\}$

$\bar{\mathcal{F}}$  is an  $\epsilon$ -net of  $\mathcal{F}$ . Also:

$$|\bar{\mathcal{F}}| = \binom{n + \lfloor L/\epsilon \rfloor}{\lfloor L/\epsilon \rfloor}$$

no. of non-negative int. solns to  $y_1 + y_2 + \dots + y_k = m$  is  $\binom{m+k-1}{k}$



Example 2 (bounded derivatives)

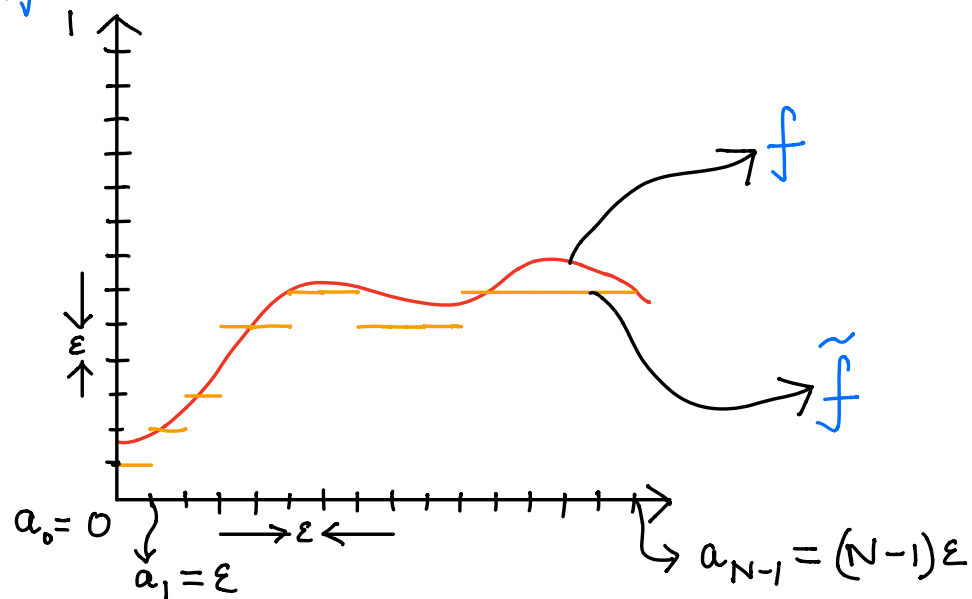
$$\mathcal{F} := \left\{ f : [0, 1] \rightarrow [0, 1], |f'| \leq 1 \right\}$$

Then, for some constant  $A < \infty$ ,

$$H_\infty(\mathcal{F}, \varepsilon) \leq A \frac{1}{\varepsilon} \quad \forall \varepsilon > 0.$$



## Proof Sketch.



$$\sup_{0 \leq x \leq 1} |f(x) - \tilde{f}(x)| \leq 2\epsilon.$$

$$|\tilde{f}(a_k) - \tilde{f}(a_{k-1})| \leq 3\epsilon, \quad k=1, 2, \dots$$

$$N_\infty(\mathcal{F}, \epsilon) \leq \left( \lfloor \frac{1}{\epsilon} \rfloor + 1 \right) 7^{\lfloor \frac{1}{\epsilon} \rfloor}$$



### Example 3 (Finite Dimensional Space)

Suppose  $\psi_1, \psi_2, \dots, \psi_d \in L_2(Q)$

and

$$\mathcal{F} := \left\{ \begin{array}{l} f = \sum_{k=1}^d \theta_k \psi_k : \theta = (\theta_1, \theta_2, \dots, \theta_d) \in \mathbb{R}^d \\ \text{and } \|f\|_{2,Q} \leq R \end{array} \right\}$$

Then,

$$H_2(\mathcal{F}, \varepsilon, Q) \leq d \log \left( \frac{4R + \varepsilon}{\varepsilon} \right)$$



### Example 3 (Finite Dimensional Space)

Suppose  $\psi_1, \psi_2, \dots, \psi_d \in L_2(Q)$

and

$$\mathcal{F} := \left\{ \begin{array}{l} f = \sum_{k=1}^d \theta_k \psi_k : \theta = (\theta_1, \theta_2, \dots, \theta_d) \in \mathbb{R}^d \\ \text{and } \|f\|_{2,Q} \leq R \end{array} \right\}$$

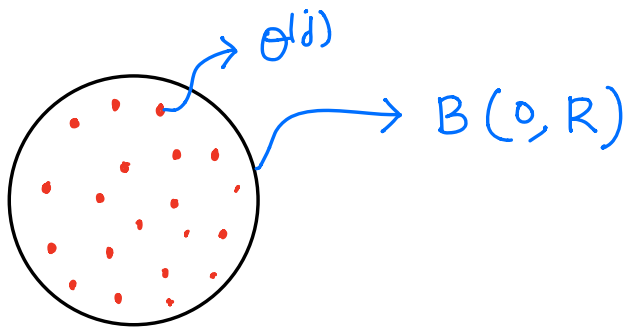
Then,

$$H_2(\mathcal{F}, \varepsilon, Q) \leq d \log \left( \frac{4R^2 + \varepsilon}{\varepsilon} \right)$$

## Proof Sketch



Find balls  $B(\theta^{(j)}, \frac{\epsilon}{R})$ ,  $j=1, 2, \dots, N$  such that the  $B(0, R)$  ball in  $\mathbb{R}^d$  with Euclidean metric is covered.



$$\left\{ \tilde{f} : \tilde{f} = \sum_{k=1}^d \theta_k^{(j)} \psi_k, j=1, 2, \dots, N \right\}$$

And,

$$N \leq \left( \frac{4R + \epsilon/R}{\epsilon/R} \right)^d.$$

Other Examples (Birman & Solomijak, 1967) (III)

$$1. \mathcal{F} := \left\{ f: \mathbb{R} \rightarrow [0, 1], \underbrace{TV(f)}_{\int_x |df|} \leq 1 \right\}$$

then,  $\exists A$  s.t.

$$H_B(\mathcal{F}, \varepsilon, Q) \leq \frac{A}{\varepsilon} \quad \forall \varepsilon > 0.$$

$$2. \mathcal{F} := \left\{ f: [0, 1] \rightarrow [0, 1], \int (f^{(m)})^2 \leq 1 \right\}$$

then,  $\exists A$  s.t.

$$H_B(\mathcal{F}, \varepsilon, Q) \leq \frac{A}{\varepsilon^{1/m}} \quad \forall \varepsilon > 0.$$



IV

## The Basic ULLN

Suppose  $\mathcal{F}$  is such that

$$H_{1,B}(\mathcal{F}, \varepsilon, P) < \infty \quad \forall \varepsilon > 0.$$

Then

$$\|P_n - P\|_\infty = \sup_{f \in \mathcal{F}} |P_n f - P f| \xrightarrow{\text{a.s.}} 0.$$

## Proof Sketch

(IV)

We can find  $\{[f_j^L, f_j^U]_{j=1}^N\}$  so that  
for each  $f \in \mathcal{F}$ , we find  $i$  s.t.

$$P_n f - P f \leq (P_n - P) f_j^U + \varepsilon$$

$$P_n f - P f \geq (P_n - P) f_j^L - \varepsilon \quad \text{--- (1)}$$

Since  $N < \infty$ , for large  $n$ ,

$$\max_{1 \leq j \leq N} |(P_n - P) f_j^U| \leq \varepsilon$$

$$\max_{1 \leq j \leq N} |(P_n - P) f_j^L| \leq \varepsilon \quad \text{--- (2)}$$

Use (1) and (2).

# AGENDA

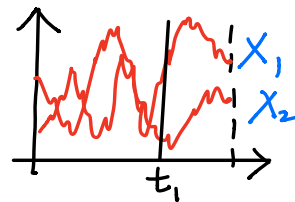
## I. Strengthen Basic ULLN.

Suppose  $\int \sup_{f \in \mathcal{F}} \|f\|_{\infty} dP < \infty$  and

$$\frac{1}{n} H_1(\mathcal{F}, \varepsilon, P_n) \xrightarrow{P} 0 \quad \forall \varepsilon > 0.$$

Then,

$$\|P_n - P\|_{\infty} := \sup_{f \in \mathcal{F}} |P_n f - P f| \xrightarrow{\text{a.s.}} 0$$



## II. Example Classes (VC)

## TWO   KEY   MACHINERY

~~(i)~~

Chaining

(ii)

Symmetrization

## SYMMETRIZATION.

means approximating

$$\|P_n - P\|_\infty \text{ using } \|P_n - P'_n\|_\infty$$

where:

$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P$$

$$X'_1, X'_2, \dots, X'_n \stackrel{\text{iid}}{\sim} P$$

$$(X_1, X_2, \dots, X_n) \stackrel{\text{ind}}{\sim} (X'_1, X'_2, \dots, X'_n)$$

## SYMMETRIZATION LEMMA

Suppose that  $\forall f \in \mathcal{F}, \delta > 0$

$$P\left(|P_n f - P f| > \delta/2\right) \leq \frac{1}{2}.$$

Then,

$$\begin{aligned} P\left(\|P_n - P\|_\infty > \delta\right) \\ \leq 2P\left(\|P_n - P'_n\|_\infty > \delta/2\right) \end{aligned}$$

## Proof Sketch.

$$P\left(\sup_{f \in \mathcal{F}} \left| \int f d(P_n - P) \right| > \delta\right)$$

$$\leq P\left(\left| \int f^* d(P_n - P) \right| > \delta\right)$$

$$\leq 2P\left(\left| \int f^* d(P_n - P) \right| > \delta \cap \left| \int f^* d(P'_n - P) \right| \leq \frac{\delta}{2}\right)$$

$$\leq 2P\left(\left| \int f^* d(P_n - P'_n) \right| > \delta/2\right)$$

$$\leq 2P\left(\sup_{f \in \mathcal{F}} \left| \int f d(P_n - P'_n) \right| > \frac{\delta}{2}\right)$$

## Why symmetrization?

$$(W_1, W_2, \dots, W_n) \stackrel{\text{ind}}{\sim} (X_1, X_2, \dots, X_n, X'_1, \dots, X'_n)$$

(Rademacher sequence)  $W_j \stackrel{\text{iid}}{\sim} \begin{cases} 1 & \text{wp } \frac{1}{2} \\ -1 & \text{wp } \frac{1}{2} \end{cases}$

Then, for each  $f \in \mathcal{F}$ ,

$$\left\{ f(X_i) - f(X'_i) : i = 1, 2, \dots, n \right\}$$

$$\stackrel{d}{=} \left\{ W_i (f(X_i) - f(X'_i)) : i = 1, 2, \dots, n \right\}$$

Okay. So what?



Notice:

$$P\left(\|P_n - P'_n\|_\infty > \frac{\delta}{2}\right)$$

$$\stackrel{\text{def.}}{=} P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - f(X'_i) \right| > \frac{\delta}{2}\right)$$

$$= P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n W_i (f(X_i) - f(X'_i)) \right| > \frac{\delta}{2}\right)$$

$$\leq 2 P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n W_i f(X_i) \right| > \frac{\delta}{4}\right).$$

The tail probability

$$P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n W_i f(X_i) \right| > \frac{\delta}{4}\right)$$

is easy to work with.

Okay, so what?

Suppose that we show:

$$P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n W_i f(X_i) \right| > \frac{\delta}{4}\right) \xrightarrow{P} 0.$$

———— (1)

Then,  $\|P_n - P\|_\infty \xrightarrow{P} 0.$   
(symmetrization)

$\|P_n - P\|_\infty \xrightarrow{\text{a.s.}} \text{something.}$   
(Pollard, 1984)

Therefore,  $\|P_n - P\|_\infty \xrightarrow{\text{a.s.}} 0.$

Let's prove (1) for bounded r.v.s.

Assume  $\sup_{f \in \mathcal{F}} \|f\|_{\infty} \leq R < \infty$ .

We will choose  $A_n \subseteq \mathcal{X}^n$  s.t.

$$P\left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n W_i f(X_i) \right| > \frac{\delta}{4}\right)$$

$$\leq P\left(E_n \mid (X_1, X_2, \dots, X_n) \in A_n\right) \times 1$$

$$+ 1 \times P\left((X_1, X_2, \dots, X_n) \in A_n^c\right)$$

$\rightarrow 0$ .

Choose :

$$A_n = \left\{ \sqrt{n} \frac{\delta}{8} \geq c \left( R H_2^{1/2} \left( \mathcal{F}, \frac{\delta}{32}, P_n \right) \sqrt{R} \right) \right\}$$

With some work, and Hoeffding (1963)

$$P(E_n \mid (X_1, X_2, \dots, X_n) \in A_n) \leq c \exp \left\{ -\frac{n\delta^2}{64c^2R^2} \right\}$$

And :

$$P(A_n^c) \rightarrow 0 \text{ if } \frac{1}{n} H_2 \left( \mathcal{F}, \delta, P_n \right) \xrightarrow{P} 0. \quad \forall \delta > 0.$$

## Key Lemma

→ "envelope"

Suppose  $F := \sup_{f \in \mathcal{F}} \|f\|_{\infty} \leq R < \infty$ ,

and that

$$\frac{1}{n} H_2(\mathcal{F}, \delta, P_n) \xrightarrow{P} 0, \quad \forall \delta > 0.$$

Then,

$$\|P_n - P\|_{\infty} \xrightarrow{\text{a.s.}} 0.$$

(7)

Theorem.

Suppose  $\int F dP < \infty$ ,

and suppose

$$\frac{1}{n} H_1(\mathcal{F}, \delta, P_n) \xrightarrow{P} 0, \quad \forall \delta > 0.$$

— (Ent)

Then,

$$\|P_n - P\|_{\infty} \xrightarrow{\text{a.s.}} 0.$$

(7)

Proof Sketch. Choose  $R$  so that

$$\|P_n - P\|_\infty := \sup_{f \in \mathcal{F}} |P_n f - P f|$$

$$\leq \sup_{f \in \mathcal{F}} \left| \int_{F \leq R} f d(P_n - P) \right|$$

$$+ \int_{F > R} F dP_n + \int_{F > R} F dP$$

$\leq 2\delta$  a.s. for large  $n$ .

$\leq \delta$  for large  $n$



## Proof Sketch. contd...

Let's deal with the first term.

Define the truncated class.

$$\mathcal{F}_R := \left\{ f \mathbb{I}_{\{f \leq R\}} : f \in \mathcal{F} \right\}$$

For  $f_1, f_2 \in \mathcal{F}$ ,

$$\int (f_1 - f_2)^2 dP_n \leq 2R \int |f_1 - f_2| dP_n$$

Proof Sketch. contd...

Under (Ent), this means

$$\frac{1}{n} H_2(\mathcal{F}_R, \delta, P_n) \xrightarrow{P} 0 \quad \forall \delta > 0.$$

Therefore, Glivenko-Cantelli  
holds on  $\mathcal{F}_R$ .

⑦

Which  $\mathcal{F}$ s satisfy the  
condition

$$\frac{1}{n} H_1(\mathcal{F}, \delta, P_n) \xrightarrow{P} 0 ?$$

Many useful classes  $\mathcal{F}$   
satisfy this property. Let's  
see one such.

# Vapnik - Chervonenkis (VC)

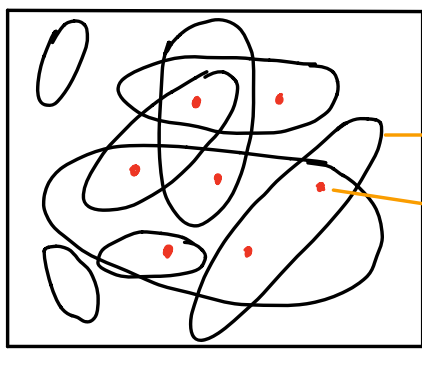
## Subgraph Classes.

$\mathcal{D} :=$  collection of subsets of  $X$

(domain of  $f \in \mathcal{F}$ )

$$\Delta^{\mathcal{D}}(X_1, X_2, \dots, X_n)$$

$$= \text{card} \left\{ D \cap \{X_1, X_2, \dots, X_n\} : D \in \mathcal{D} \right\}$$



$$\Delta^{\mathcal{D}}(X_1, X_2, \dots, X_7) = 6.$$

$$m^{\mathcal{D}}(n) :=$$

$$\sup \left\{ \Delta^{\mathcal{D}}(x_1, x_2, \dots, x_n) : (x_1, x_2, \dots, x_n) \in \mathcal{X} \right\}$$

$$V(\mathcal{D}) := \inf \left\{ n \geq 1 : m^{\mathcal{D}}(n) < 2^{-n} \right\}$$

-  $V(\mathcal{D})$  is called the index of the class  $\mathcal{D}$ .

-  $\mathcal{D}$  is called a **VC-class** if  $V(\mathcal{D}) < \infty$ .

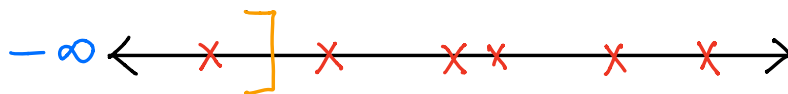
## Example 1

$$\mathcal{D} := \{(-\infty, t], t \in \mathbb{R}\}$$

Then,

$$m^{\mathcal{D}}(n) = n+1, \text{ so that}$$

$\mathcal{D}$  is a VC-class.



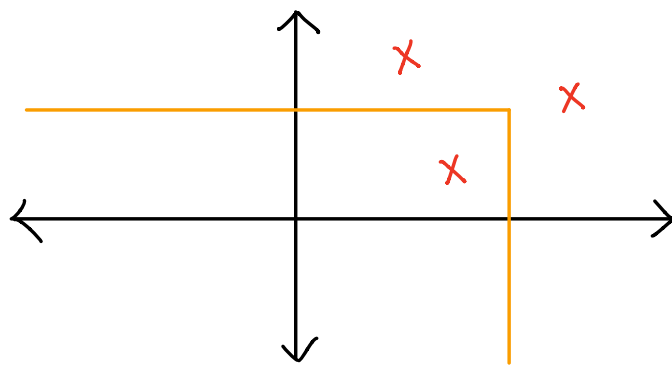
## Example 2

$$\mathcal{D} := \left\{ (-\infty, t], t \in \mathbb{R}^d \right\}$$

Then,

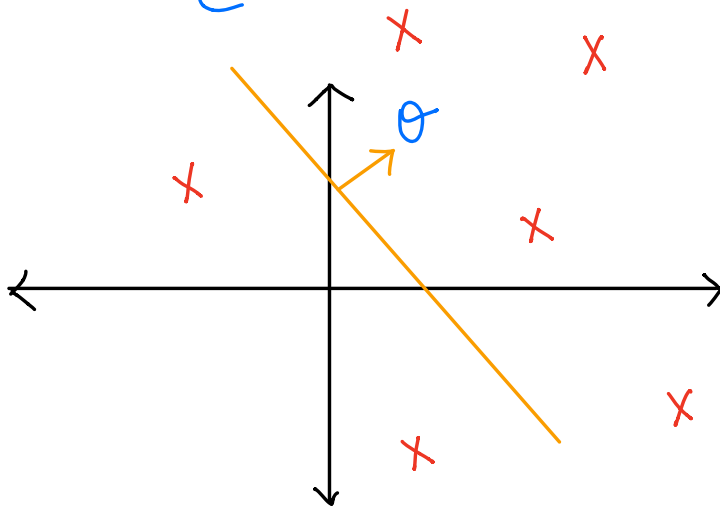
$$m^{\mathcal{D}}(n) = (n+1)^d, \text{ so that}$$

$\mathcal{D}$  is a VC-class.



### Example 3

$$\mathcal{D} := \left\{ \left\{ x : \langle x, \theta \rangle > y \right\} : \theta \in \mathbb{R}^d, y \in \mathbb{R} \right\}.$$



$$m^{\mathcal{D}}(n) \leq 2^d \binom{n}{d}$$

$$V(\mathcal{D}) \leq d+2.$$