

Proof.  $x^{(g)}(f) \in C_{\text{lf}}^{\infty}(\mathbb{R})$

$$|x^{(g)}(f)| \leq C \int_{-\infty}^{\infty} |x'(f)| |f|^g \leq C \int_{-\infty}^{\infty} |\zeta(f)| (1+|f|)^g df$$

$$\Rightarrow \int_{-\infty}^{\infty} |\zeta(f)| (1+|f|)^g df < \infty \text{, for every } g \leq p.$$

$$|\zeta'(f)| \leq \frac{C'}{(1+|f|)^{p+1+\varepsilon}}$$

### III.1.3. Non-linear Approximation

B. - ONB

$$x = \sum_{j=0}^{\infty} \langle x, g_j \rangle g_j \quad | \quad x_n = \sum_{j \in \mathbb{N}_n} \langle x, g_j \rangle g_j$$

$$\begin{aligned} \mathcal{E}[n] &= \|x - x_n\|^2 = \left\| \sum_{j=0}^{\infty} \langle x, g_j \rangle g_j - \sum_{j \in \mathbb{N}_n} \langle x, g_j \rangle g_j \right\|^2 \\ &= \left\| \sum_{j \notin \mathbb{N}_n} \langle x, g_j \rangle g_j \right\|^2 = \sum_{j \notin \mathbb{N}_n} |\langle x, g_j \rangle|^2 \end{aligned}$$

$$\|x\|^2 = \sum_j |\langle x, g_j \rangle|^2 = \sum_{j \in \mathbb{N}_n} |\langle x, g_j \rangle|^2 + \sum_{j \notin \mathbb{N}_n} |\langle x, g_j \rangle|^2$$

$$\mathcal{E}[n] = \|x\|^2 - \sum_{j \in \mathbb{N}_n} |\langle x, g_j \rangle|^2$$

sort  $\{|\langle x, g_j \rangle|\}_{j \in \mathbb{N}}$  in descending order and denote  $x_B^{[k]} = \langle x, g_k \rangle$

$$|x_B^{[k]}| \geq |x_B^{[k+1]}|, k \geq 1$$

$$x_n = \sum_{k=1}^n x_B^{[k]} g_k$$

$$\mathcal{E}[n] = \|x - x_n\|^2 = \sum_{k=n+1}^{\infty} |x_B^{[k]}|^2$$

Thm 11.4. Let  $s > 1/2$ . If there exists  $C > 0$  s.t.  $|x_B^r(s)| \leq C s^{-s}$ , then

$$E_n(M) \leq \frac{C^2}{2^{s-1}} M^{1-s/2}. \quad (*)$$

Conversely, if  $E_n(M)$  satisfies  $(*)$ , then

$$|x_B^r(s)| \leq (1 - \frac{1}{2^s})^{-s} C s^{-s}. \quad (**)$$

$$\|x\|_{B,p} = \left( \sum_{j=0}^{\infty} |(x, g_j)|^p \right)^{1/p}$$

Thm 11.5. Let  $p < 2$ . If  $\|x\|_{B,p} < \infty$ , then

$$|x_B^r(s)| \leq \|x\|_{B,p} s^{-1/p}$$

and

$$E_n(M) = o(M^{1-2/p}), \text{ i.e.,}$$

$$\text{i.e. } \lim_{M \rightarrow \infty} E_n(M) M^{-1+2/p} = 0.$$

$$B_{B,p} = \{x \in H : \|x\|_{B,p} < \infty\}$$

Summary. If  $x \in B_{B,p}$ , then  $E_n(M) = o(M^{1-2/p})$ . This is called a "Jackson type inequality". Conversely, if  $E_n(M) = o(M^{1-2/p})$ ,

then the "Bernstein inequality" guarantees  $|x_B^r(s)| \leq C s^{-1/p}$ , which shows that  $x \in B_{B,q}$  for every  $q > p$ , we have

$$\sum_s |x_B^r(s)|^q \leq \sum_s C^q s^{-q/p} < \infty.$$

E.g. for wavelike bases, the space  $B_{B,p}$  is a Besov space.

WH bases, the space  $B_{B,p}$  is a Feichtinger-Fröbenius modulor space.

## Chapter 12. Uniform Laws of Large numbers

### 12.1. Uniform convergence of CDFs

$X \dots \text{RV}$

$$\bar{F}(t) = P[X \leq t], t \in \mathbb{R}.$$

Problem. We want to estimate  $\bar{F}(t)$  from a given collection  $\{x_i\}_{i=1}^n$  of i.i.d. samples each drawn according to the law specified by  $\bar{F}$ .

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(-\infty, t]}(x_i), \text{ empirical CDF}$$

where  $\mathbf{1}_{(-\infty, t]}(x)$  is the indicator function of the event  $\{x \leq t\}$ .

$$\begin{aligned} \bar{F}(t) &= E[\mathbf{1}_{(-\infty, t]}(X)] \leftarrow \text{population CDF} \\ &= \int_{-\infty}^t \mathbf{1}_{(-\infty, t]}(x) f(x) dx \\ &= \int_{-\infty}^t f(x) dx \quad (= \bar{F}(t)) \end{aligned}$$

want to understand the properties of  $\hat{F}_n(t)$ .

$$\begin{aligned} E[\hat{F}_n(t)] &= \frac{1}{n} \sum_{i=1}^n E[\mathbf{1}_{(-\infty, t]}(X_i)] \\ &= \frac{1}{n} \sum_{i=1}^n \bar{F}(t) = \bar{F}(t) \end{aligned}$$

One can show that, indeed,  $\hat{F}_n(t)$  converges almost surely to  $\bar{F}(t)$ , i.e.,

$$P\left[\lim_{n \rightarrow \infty} \hat{F}_n(t) = \bar{F}(t)\right] = 1, \quad \forall t \in \mathbb{R}.$$

functionals  $p$  that map a CDF  $\bar{F}$  to a real number, namely  $p(\bar{F})$ , i.e.,  $\bar{F} \mapsto p(\bar{F})$

Problem: We want to estimate  $p(\bar{F})$  from  $\{x_i\}_{i=1}^n$

Q: Can we get good estimates of  $p(\bar{F})$  by replacing  $\bar{F}$  by  $\hat{F}_n$ , i.e.,

what can we say about  $p(\hat{F}_n)$ ? In particular, what can we say about the convergence of  $p(\hat{F}_n)$  to  $p(\hat{F})$  as  $n \rightarrow \infty$ .

Example 12.1. (Expectation functional). Given an integrable  $g$ , we define

$$pg(F) = E(g(X)) = \int g(x)f(x)dx = \int g(x)d\bar{F}(x).$$

$$\text{e.g. } g(x) = x \Rightarrow pg(\hat{F}) = \int x f(x)dx = E(X).$$

$pg(\hat{F}_n)$  relative to  $pg(\hat{F})$



plug-in estimator

$$pg(\hat{F}_n) = \frac{1}{n} \sum_{i=1}^n g(x_i)$$

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, x]}(x_i)$$

$$\frac{d\bar{F}(x)}{dx} = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$$

↑  
Dirac delta

$$\mathbb{1}_{(-\infty, x]}(x_i) = \begin{cases} 1, & x_i \leq x \\ 0, & \text{else} \end{cases}$$



plug-in  $\hat{F} \rightarrow \hat{F}_n$

$$\int g(x)d\bar{F}(x) = \frac{1}{n} \sum_{i=1}^n g(x_i)$$

↑  
to be shown

$$\int g(x) \frac{1}{n} \sum_{i=1}^n \delta(x - x_i) dx = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} g(x) \delta(x - x_i) dx$$

$$= \frac{1}{n} \sum_{i=1}^n g(x_i).$$

Example 12.2. (Quantile functions). For any  $\alpha \in (0, 1)$ , the quantile functional  $Q_\alpha$  is given by

$$Q_2(F) = \inf\{t \in \mathbb{R} \mid F(t) \geq \frac{1}{2}\}.$$

(median corr. to  $\lambda=0.5$ ). The plug-in est. is

$$\hat{Q}_2(\hat{F}_n) = \inf\{t \in \mathbb{R} \mid \frac{1}{n} \sum_{i=1}^n 1_{(-\infty, t]}(X_i) \geq \frac{1}{2}\}.$$

Want to understand convergence of  $p_g(\hat{F}_n)$  to  $p_g(F)$  as  $n \rightarrow \infty$ .

Def. 12.4. We say that a sequence of random variables  $X_n$  converges in probability to the random variable  $X$  if

$$\lim_{n \rightarrow \infty} P\{|X - X_n| > \varepsilon\} = 0, \forall \varepsilon > 0.$$

a.s.  $\Rightarrow$  p. (Fatou's lemma)

$$\|G - \bar{F}\|_\infty = \sup_{t \in \mathbb{R}} |G(t) - \bar{F}(t)|, \quad \text{distance between } \bar{F} \text{ and } G$$

$$(||X||_\infty = \sup_{t \in \mathbb{R}} |X(t)|)$$

Def. 12.5. The functional  $p$  is continuous at  $\bar{F}$  in the sup-norm, if, for all  $\varepsilon > 0$ , there exists a  $\delta > 0$  s.t.  $\|G - \bar{F}\|_\infty \leq \delta \Rightarrow |p(G) - p(\bar{F})| < \varepsilon$ .



We establish that if  $p_g$  is continuous in the sup-norm (at  $\bar{F}$ ), then  $p_g(\hat{F}_n) \xrightarrow{\text{P.}} p_g(\bar{F})$ .

Proof follows by combining continuity of the functional with

Thm. 12.6. (Glivenko-Cantelli). For any distribution, the empirical CDF  $\hat{F}_n$  satisfies

$$\|\hat{F}_n - \bar{F}\|_\infty \xrightarrow{\text{a.s.}} 0.$$

We need to show that

$$\lim_{n \rightarrow \infty} P[|p_g(\hat{F}_n) - p_g(\bar{F})| > \varepsilon] = 0, \quad \forall \varepsilon > 0.$$

(conv. in prob.)

$$\|\hat{f}_n - f\|_{\infty} \leq \delta \Rightarrow |p(\hat{f}_n) - p(f)| \leq \varepsilon \quad (\text{continuity})$$

Hence, it follows that

$$|p(\hat{f}_n) - p(f)| > \varepsilon \Rightarrow \|\hat{f}_n - f\|_{\infty} > \delta$$

$\Rightarrow$

$$P[|p(\hat{f}_n) - p(f)| > \varepsilon] \leq P[\|\hat{f}_n - f\|_{\infty} > \delta]$$

Now, Glivenko - Cantelli  $\|\hat{f}_n - f\|_{\infty} \xrightarrow{\text{a.s.}} 0 \Rightarrow \|\hat{f}_n - f\|_{\infty} \xrightarrow{p.} 0$ , i.e.,

$$\lim_{n \rightarrow \infty} P[\|\hat{f}_n - f\|_{\infty} > \varepsilon] = 0, \forall \varepsilon' > 0$$

This establishes that

$$\lim_{n \rightarrow \infty} P[|p(\hat{f}_n) - p(f)| > \varepsilon] \leq \lim_{n \rightarrow \infty} P[\|\hat{f}_n - f\|_{\infty} > \delta] = 0$$

$$\Rightarrow \lim_{n \rightarrow \infty} P[|p(\hat{f}_n) - p(f)| > \varepsilon] = 0, \forall \varepsilon > 0$$

## 12.2. Uniform Laws for more general function classes

Let  $\mathcal{F}$  be a class of integrable functions with domain  $X$ , and let  $\{X_i\}_{i=1}^n$  be a collection of i.i.d. samples from some distribution  $P$  over  $X$ .

$$\|\hat{P}_n - P\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - E[f(x)] \right|$$

Def. 12.7. We say that  $\mathcal{F}$  is a Glivenko - Cantelli class for  $P$  if  $\|\hat{P}_n - P\|_{\mathcal{F}}$  converges to zero in prob. as  $n \rightarrow \infty$ .

Ex. 12.8. (Empirical CDF and indicator functions).

$$\mathcal{F} = \{1_{(-\infty, t]}(\cdot) | t \in \mathbb{R}\}.$$

For fixed  $t \in \mathbb{R}$ , we have

$$E[\mathbb{1}_{(-\infty, +)}(x)] = P[X \leq +] = \bar{F}(+)$$

GC.  $\|\hat{f}_n - f\|_{\infty} \xrightarrow{\text{a.s.}} 0$ , i.e., (a.s.  $P(\lim_{n \rightarrow \infty} x_n = x) = 1$ )

$$P\left[\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} |\hat{f}_n(t) - f(t)| = 0\right] = 1$$

↓

$$P\left[\lim_{n \rightarrow \infty} \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, +)}(x_i) - E[\mathbb{1}_{(-\infty, +)}(X)] \right| = 0\right] = 1$$

$$\boxed{= P\left[\lim_{n \rightarrow \infty} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - E[f(X)] \right| = 0\right] = 1.}$$

$$\|\hat{P}_n - P\|$$

Where does this occur?

family of prob. distributions  $\{P_\theta | \theta \in \mathcal{S}\}$

given  $\{x_i\}_{i=1}^n$  drawn i.i.d. according to  $P_{\theta^*}$ , for some fixed,  
but unknown  $\theta^* \in \mathcal{S}$

$\mathcal{S} = \mathbb{R}^d$  vector estimates (finite-dimensional) parametric

$\mathcal{S} = \mathcal{G}$  function class non-parametric

A standard approach for estimating  $\theta^*$  from  $\{x_i\}_{i=1}^n$  is based on minimizing a cost function  $L_\theta(x)$ , which quantifies the fit between  $\theta \in \mathcal{S}$  and the sample  $\{x_i\}_{i=1}^n$ .

$$\hat{R}_n(\theta, \theta^*) = \frac{1}{n} \sum_{i=1}^n L_\theta(x_i)$$

empirical risk

$$R(\theta, \theta^*) = E_{\theta^*}[L_\theta(x)]$$

population risk

in practice, typically minimize over a subset  $\mathcal{S}_0$  of  $\mathcal{S}$ ,  $\Rightarrow$  est.  $\hat{\theta}$

$$E(\hat{\theta}, \theta) = R(\hat{\theta}, \theta^*) - \inf_{\theta \in \mathcal{S}_0} R(\theta, \theta^*)$$

Ex. 12.9. (Maximum likelihood).  $\{P_\theta(\theta \in \mathcal{D})\}$ ,  $\theta^*$  unknown parameter

$$\{x_i\}_{i=1}^n$$

$$L_\theta(x) = \log \left( \frac{P_{\theta^*}(x)}{P_\theta(x)} \right)$$

ML estimator

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{D}_0} \left\{ \frac{1}{n} \sum_{i=1}^n \log \left[ \frac{P_{\theta^*}(x_i)}{P_\theta(x_i)} \right] \right\}$$

$$= \arg \min_{\theta \in \mathcal{D}_0} \left\{ \frac{1}{n} \sum_{i=1}^n \log \left[ \frac{1}{P_\theta(x_i)} \right] \right\}$$

$$\text{Population risk: } R(\theta, \theta^*) = E_{\theta^*} \left[ \log \left( \frac{P_{\theta^*}(x)}{P_\theta(x)} \right) \right] = \underbrace{D(P_{\theta^*} || P_\theta)}$$

Kullback-Leibler divergence

$$\inf_{\theta \in \mathcal{D}_0} R(\theta, \theta^*) = \inf_{\theta \in \mathcal{D}_0} E_{\theta^*} \left[ \log \left( \frac{P_{\theta^*}(x)}{P_\theta(x)} \right) \right].$$

$$D(p||q) \geq 0 \text{ with } = \text{ iff } p(x) = q(x), \forall x.$$

$$E(\hat{\theta}, \theta^*) = R(\hat{\theta}, \theta^*) - \underbrace{\inf_{\theta \in \mathcal{D}_0} R(\theta, \theta^*)}_{= 0 \text{ if } \theta^* \in \mathcal{D}_0} = R(\hat{\theta}, \theta^*) = D(P_{\theta^*} || P_{\hat{\theta}})$$

□

How does ERM relate to  $\|P_n - P\|_F$

$$\text{assume that } \hat{\theta}_0 \in \mathcal{D}_0 \text{ s.t. } R(\hat{\theta}_0, \theta^*) = \inf_{\theta \in \mathcal{D}_0} R(\theta, \theta^*)$$

$$E(\hat{\theta}, \theta^*) = R(\hat{\theta}, \theta^*) - \underbrace{R_n(\hat{\theta}, \theta^*)}_{\bar{T}_1} + \underbrace{R_n(\hat{\theta}, \theta^*) - R_n(\theta_0, \theta^*)}_{\bar{T}_2} + \underbrace{R_n(\theta_0, \theta^*) - R(\theta_0, \theta^*)}_{\bar{T}_3}$$

$\bar{T}_2 \leq 0$  because  $\hat{\theta}$  minimizes ER over  $\mathcal{D}_0$ .

$$\bar{T}_3 = \frac{1}{n} \sum_{i=1}^n L_{\theta_0}(x_i) - E_x [L_{\theta_0}(x)] \quad \text{concentration of measure meas.}$$

$$\bar{T}_1 = E_x [L_{\theta}(x)] - \frac{1}{n} \sum_{i=1}^n L_{\theta}(x_i)$$

challenging because  $\hat{\theta}$  is a random quantity

$$\mathcal{G} = \{L_{\theta}(\cdot) | \theta \in \Theta_0\}$$

$$|\bar{T}_1| \leq \sup_{\theta \in \Theta_0} \left| \frac{1}{n} \sum_{i=1}^n L_{\theta}(x_i) - E_x [L_{\theta}(x)] \right| = \|P_n - P\|_{\mathcal{G}}$$

$$E(\bar{\theta}_1 | \theta^*) \leq 2 \|P_n - P\|_{\mathcal{G}}.$$