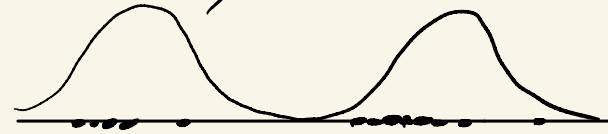


# GLIVENKO - CANTELLI THEOREM

• Basic problem in statistics: what can we learn about the unknown distribution from a random sample?

• Mean: Let  $X, X_1, X_2, \dots$  be iid r.v.s with finite mean  $\mu$ .

$$\text{SLLN: } \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{a.s.}} \mu.$$



• Entire distribution? i.e. is

$$F(x) = P\{X \leq x\} \quad (\text{cdf})$$

$$\stackrel{?}{\approx} F_n(x) = \frac{1}{n} \#\{i=1, \dots, n: X_i \leq x\} \quad (\text{empirical cdf}) \quad ?$$

$$\forall x \in \mathbb{R}: F_n(x) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{1}_{\{X_i \leq x\}}}_{\substack{\uparrow \\ \text{iid, mean} = P\{X \leq x\}}} \xrightarrow[\text{SLLN}]{\text{a.s.}} P\{X \leq x\} = F(x)$$

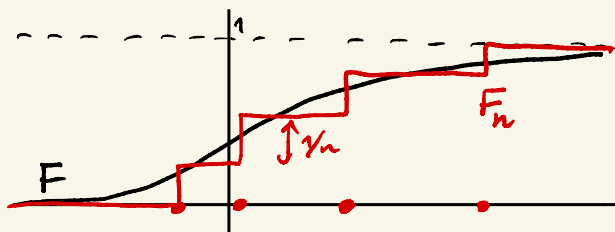
• Improvement from pointwise to uniform convergence:

Thm [Glivenko-Cantelli '1933] let  $X_1, X_2, \dots$  be iid r.v.s. Then

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

• GCT allows the statistician to interact with data:

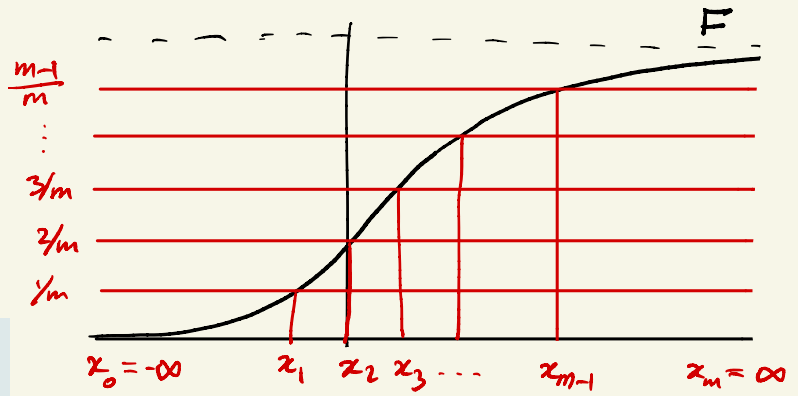
estimate  $P\{a \leq X \leq b\}$  where  $a, b$  depend on data.



Proof ① For continuous F:

Fix  $m \in \mathbb{N}$ , choose  $x_j \in \mathbb{R} \cup \{\pm\infty\}$ :

$$F(x_j) = \frac{j}{m}, \quad j=0, \dots, m.$$



$x_0$  can be  $-\infty$ ,  $x_m$  can be  $+\infty$   
 $F(-\infty) = 0$ ,  $F(+\infty) = 1$   
 well defined since  $F$  is continuous and  
 $F(x) \rightarrow 0$  as  $x \rightarrow -\infty$ ,  $F(x) \rightarrow 1$  as  $x \rightarrow \infty$

Then

$$F(x_j) - F(x_{j-1}) = \frac{1}{m}, \quad j=1, \dots, m. \quad (*)$$

• SLLN  $\Rightarrow \forall j, |F_n(x_j) - F(x_j)| \xrightarrow{a.s.} 0$  (see p. 100)  $\Rightarrow$

$$\varepsilon_n := \max_j |F_n(x_j) - F(x_j)| \xrightarrow{a.s.} 0 \quad (**)$$

•  $\forall x \in \mathbb{R} \exists j: x_{j-1} \leq x \leq x_j \Rightarrow$

$$F_n(x) - F(x) \leq F_n(x_j) - F(x_{j-1}) \quad (\text{monotonicity})$$

$$\leq F_n(x_j) - F(x_j) + \frac{1}{m} \quad (\text{by } *)$$

$$\leq \varepsilon_n + \frac{1}{m} \quad (\text{by } **)$$

Similarly,

$$F_n(x) - F(x) \geq F_n(x_{j-1}) - F(x_j) \quad (\text{monotonicity})$$

$$\geq F_n(x_{j-1}) - F(x_{j-1}) + \frac{1}{m} \quad (\text{by } *)$$

$$\geq -\varepsilon_n + \frac{1}{m} \quad (\text{by } **).$$

$\Rightarrow$

$$|F_n(x) - F(x)| \leq \varepsilon_n + \frac{1}{m} \quad \forall x \in \mathbb{R} \Rightarrow$$

$$\|F_n - F\|_\infty \leq \varepsilon_n + \frac{1}{m} \quad \forall n \in \mathbb{N}$$

(\*\*)  $\Rightarrow$  with prob. 1,  $\limsup_n \|F_n - F\|_\infty \leq \frac{1}{m} \quad \forall m \in \mathbb{N}$

with prob. 1,  $\|F_n - F\|_\infty \rightarrow 0$ .  $\square$

② For arbitrary  $F$ ,

Use right continuity:

$$F(x^+) := \lim_{y \downarrow x} F(y) = F(x)$$

Also,

$$F(x^-) := \lim_{y \uparrow x} F(y) = \lim_{y \uparrow x} P\{X \leq y\} = P\{X < x\} \quad (\text{by continuity of probability})$$

Fix  $m \in \mathbb{N}$ , define

$$x_j := \inf \left\{ x : F(x) \geq \frac{j}{m} \right\}, \quad j=0, \dots, m$$

(as before,  $x_0$  can be  $-\infty$ ,  $x_m$  can be  $+\infty$ )

• Claim:  $F(x_j^-) - F(x_{j-1}) \leq \frac{1}{m}, \quad j=1, \dots, m. \quad (*)$

$$\left. \begin{aligned} F(x_j^-) &= \lim_{x \uparrow x_j} F(x) \leq \frac{j}{m} \\ &\quad \text{\small \(\wedge\) def of } x_j \\ &\quad \text{\small } j/m \\ F(x_{j-1}) &= \lim_{x \downarrow x_{j-1}} F(x) \geq \frac{j-1}{m} \\ &\quad \text{\small \(\vee\) def of } x_{j-1} \\ &\quad \text{\small } \frac{j-1}{m} \end{aligned} \right\} \text{Moreover, by right continuity,} \\ \Rightarrow F(x_j^-) - F(x_{j-1}) \geq \frac{j}{m} - \frac{j-1}{m} = \frac{1}{m}. \quad \square$$

• SLLN  $\Rightarrow$

$$\begin{aligned} \varepsilon_n &:= \max_j |F_n(x_j) - F(x_j)| \xrightarrow{\text{a.s.}} 0 \quad \leftarrow (\text{as before}) \\ \delta_n &:= \max_j |F_n(x_j^-) - F(x_j^-)| \xrightarrow{\text{a.s.}} 0 \end{aligned} \quad (**)$$

$$\left( \lim_{y \uparrow x} \frac{1}{n} \sum_{i=1}^n P\{X_i \leq y\} = \frac{1}{n} \sum_{i=1}^n P\{X_i < x\} \xrightarrow{\text{SLLN}} P\{X < x\} = F(x^-) \right)$$

•  $\forall x \in \mathbb{R} \exists j: x_{j-1} \leq x \leq x_j \Rightarrow$

$$\begin{aligned} F_n(x) - F(x) &\leq F_n(x_j^-) - F(x_{j-1}) \quad (\text{monotonicity}) \\ &\leq F_n(x_j^-) - F(x_j^-) + \frac{1}{m} \quad (\text{by } *) \\ &\leq \delta_n + \frac{1}{m} \quad (\text{by } **) \end{aligned}$$

Similarly,

$$\begin{aligned} F_n(x) - F(x) &\geq F_n(x_{j-1}) - F(x_j^-) \quad (\text{monotonicity}) \\ &\geq F_n(x_{j-1}) - F(x_{j-1}) - \frac{1}{m} \quad (\text{by } *) \\ &\geq -\varepsilon_n - \frac{1}{m} \quad (\text{by } **) \end{aligned}$$

•  $\Rightarrow |F_n(x) - F(x)| \leq \max(\varepsilon_n, \delta_n) + \frac{1}{m}$ . Finish the proof as in ①.  $\square$

## Remarks:

1.  $F_n(x) = \frac{1}{n} \#\{i \leq n: X_i \leq x\}$  is the cdf of the prob. measure

$$\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

This random meas.  $\mu_n$  is called the empirical distribution.

2. GCT can be stated as follows:

$$\sup_B \left| \frac{1}{n} \#\{i \leq n: X_i \in B\} - P\{X \in B\} \right| \xrightarrow{\text{a.s.}} 0$$

where the  $\sup$  is over all intervals  $B = (-\infty, x]$ .

$\Rightarrow$  same is true for intervals  $B = (a, b]$  or  $[a, b]$  (Why?)

Question: is this true for all Borel sets  $B$ ?

No:  $B = \{X_1, \dots, X_n\}$   $|1-0| \not\rightarrow 0$

What families of sets  $B$  is this true for?

Ans:  $\forall$  family with finite VC dimension (see e.g. my book)

$\downarrow$   
ML

3. Rate of convergence in GCT is  $O(1/\sqrt{n})$ :

Dvoretzky-Kiefer-Wolfowitz inequality '1956:

$$P\{\|F_n - F\|_\infty > \varepsilon\} \leq 2e^{-n\varepsilon^2}$$

$\uparrow$   
[Massart 1990]

4.  $\Rightarrow$  Kolmogorov-Smirnov test for whether 2 samples

come from the same distribution (e.g. if a drug has any effect):

(in this case,  $\|F_n - G_n\|_\infty \leq \|F_n - F\|_\infty + \underbrace{\|F - G\|_\infty}_b + \|G_n - G\|_\infty = O(1/\sqrt{n})$ )