

LECTURE 6

Previous class: Hoeffding's inequality:

If X_1, \dots, X_N are independent r.v.'s such that $X_i \in [a_i, b_i] \forall i$,
 then $S_N = \sum_{i=1}^N X_i$ satisfies

$$P\{|S_N - ES_N| \geq t\} \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^N (b_i - a_i)^2}\right) \quad \forall t \geq 0$$

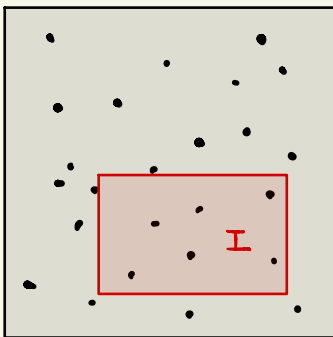
• Example: $X_i \sim \text{Ber}(p) \Rightarrow a_i = 0, b_i = 1; ES_N = pN, t := \delta N$

Hence $S_N \sim \text{Binom}(N, p)$ satisfies

$$P\{|S_N - pN| \geq \delta N\} \leq 2 \exp(-2\delta^2 N)$$

↪ exponentially small in N . 😊

Today: an application for: DISCREPANCY

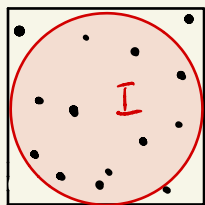


- Throw N random points into the square $[0,1]^2$ independently and uniformly.
- \forall subset $I \subset [0,1]^2$, expected fraction of pts in $I = \text{area}(I)$
- Why? $N_I := \#(\text{pts in } I) \sim \text{Binom}(N, P_I)$

$$P_I = P\{\text{a random uniform pt} \in I\} = \text{area}(I)$$

$$\Rightarrow EN_I = P_I N \quad \Rightarrow \quad E\left[\frac{N_I}{N}\right] = P_I = \text{area}(I)$$

• Application: a probabilistic computation of π :



$$\frac{N_I}{N} \xrightarrow{\text{LLN}} E\left[\frac{N_I}{N}\right] = \text{area}(I) = \frac{\pi}{4}$$

• Mean squared error: $E\left(\frac{N_I}{N} - P_I\right)^2 = \text{Var}\left(\frac{N_I}{N}\right) = \frac{1}{N^2} \text{Var}(N_I) = \frac{NP(1-P)}{N^2} \leq \frac{1}{N}$

\Rightarrow st dev $\leq \frac{1}{\sqrt{N}}$. Chebyshev \Rightarrow

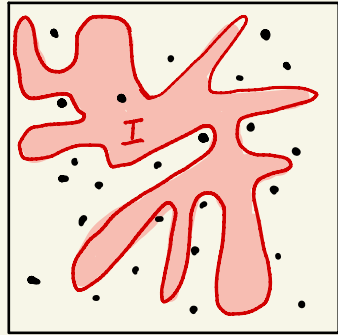
$$P\left\{\left|\frac{N_I}{N} - P_I\right| = O\left(\frac{1}{\sqrt{N}}\right)\right\} \geq 0.99 \quad \forall I \subset [0,1]^2 \quad (*)$$

• Q: does (*) hold for all I simultaneously, i.e. is it true that

$$P \left\{ \forall I \subset [0,1]^2 : \left| \frac{N_I}{N} - \text{Area}(I) \right| = O(1/\sqrt{N}) \right\} \geq 0.99? \quad (**)$$

(is there a "universal sample"?)

• No:



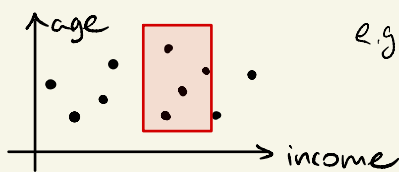
• Q: Does this hold for simple shapes, such as $I \in$ rectangles? YES 😊

THM (Discrepancy) A set of N independent random points, uniformly drawn from the square $[0,1]^2$, satisfies the following with probability ≥ 0.99 .

For any axis-aligned rectangle $I \subset [0,1]^2$, the fraction of the points in I satisfies

$$\left| \frac{N_I}{N} - \text{area}(I) \right| \leq C \sqrt{\frac{\log N}{N}}$$

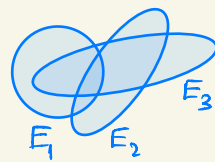
• Relevance for statistic: representative sampling:



e.g. We want a sample in which all age brackets and income brackets we fairly represented

PROOF "An epsilon-net method," based on a union bound

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) \leq \sum_{i=1}^{\infty} P(E_i) \quad \forall \text{ events } E_i$$



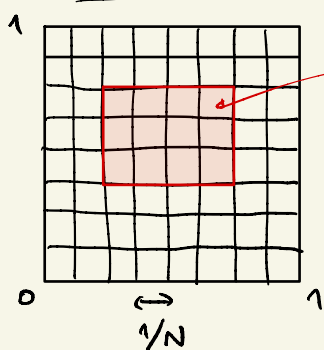
Want to show:

$$P\left\{ \exists \text{ rectangle } I \quad \left| \frac{N_I}{N} - p_I \right| > C \sqrt{\frac{\log N}{N}} \right\} \leq 0.01$$

$$= P\left(\bigcup_{I \in \text{rectangles}} E_I \right) \leq \sum_{I \in \text{rectangles}} P(E_I)$$

↑ infinite sum ☹️ ⇒

① Discretize:



a "grid rectangle"

There are $\leq N^4$ grid rectangles. Not anymore! 😊

② Concentration: \forall fixed grid rectangle I , $N_I \sim \text{Binom}(N, p_I) \Rightarrow$

$$P\left\{ \left| \frac{N_I}{N} - p_I \right| \geq \delta \right\} = P\left\{ |N_I - p_I N| \geq \delta N \right\}$$

Hoeffding p.2

$$\leq 2 \exp(-2\delta^2 N) \quad \forall \delta \geq 0$$

$$\leq \frac{1}{100 N^4}$$

if we choose

$$\delta = C \sqrt{\frac{\log N}{N}}$$

↑ a large absolute constant

③ Union bound:

union bd

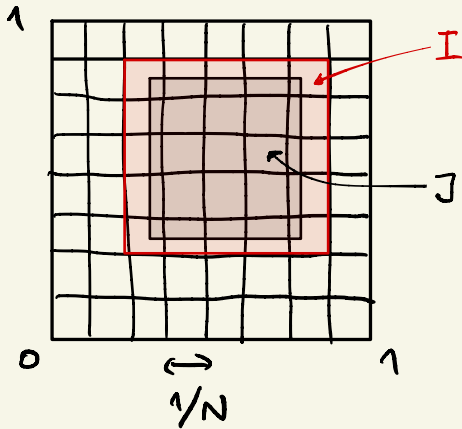
$$P\left\{ \exists I \in \text{grid rectangles} : \left| \frac{N_I}{N} - p_I \right| \geq \delta \right\} \leq \sum_{I \in \text{grid rec}} P\left\{ \left| \frac{N_I}{N} - p_I \right| \geq \delta \right\}$$

$$\leq N^4 \cdot \frac{1}{100 N^4} = 0.01.$$

We proved: $P\left\{ \forall \text{ grid rec } I : \left| \frac{N_I}{N} - p_I \right| \leq C \sqrt{\frac{\log N}{N}} \right\} \geq 0.99.$ 😊

Assume this event occurs ↗

④ Approximation " \forall rectangle \approx a grid rectangle " :



• \forall rectangle J lies in a grid rectangle I with area $P_I \leq P_J + \frac{4}{N} \Rightarrow$

$$\frac{N_J}{N} \leq \frac{N_I}{N} \leq P_I + C \sqrt{\frac{\log N}{N}} \quad (\text{by step 3})$$

$$\leq P_J + \frac{4}{N} + C \sqrt{\frac{\log N}{N}} \leq P_J + C' \sqrt{\frac{\log N}{N}}$$

↑ smaller ↑ larger

• Similarly, $\frac{N_J}{N} \geq P_J - C' \sqrt{\frac{\log N}{N}} \quad (\text{DIY})$

Thus:

$$\left| \frac{N_J}{N} - P_J \right| \leq C' \sqrt{\frac{\log N}{N}} \quad \forall \text{ rectangle } I. \quad \text{QED.}$$

REMARKS

- ① $\log N$ can be removed.
- ② Uniform distr. on $[0,1]^2$ can be replaced with \forall distr. on \mathbb{R}^2
- ③ Rectangles can be replaced by other simple shapes such as triangles, circles, ellipses...
- ④ The result can be extended to \mathbb{R}^d :

$$\left| \frac{N_I}{N} - P_I \right| \leq C \sqrt{\frac{d}{N}} \quad \forall \text{ box } I$$

All these will follow from general VC theory (covered later).