

# Lecture 12. Concentration of measure.

\* Recap lecture 11.

Claim. ERM is agnostic PAC learner with

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log |\mathcal{H}|}{2\epsilon^2 \delta} \right\rceil$$

Corollary 1:  $(\frac{\epsilon}{2}, \delta)$ -Uniform convergence  $\rightarrow$  PAC

$$\iff m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{u.c.}(\frac{\epsilon}{2}, \delta)$$

left to show:  $m \geq \left\lceil \frac{\log |\mathcal{H}|}{2\epsilon^2 \delta} \right\rceil \rightarrow$  Uniform convergence w.p  $1-\delta$

$$\mathbb{P}_S \left( \bigcup_{h \in \mathcal{H}} \{ |L_S(h) - L_D(h)| > \epsilon \} \right) \leq \delta$$

n. t. s.  $|\mathcal{H}| \cdot \mathbb{P}_S ( |L_S(h) - L_D(h)| > \epsilon ) \leq \delta$

$$m \geq \left\lceil \frac{\log |\mathcal{H}|}{2\epsilon^2 \delta} \right\rceil \text{ n. t. s. } \rightarrow \mathbb{P}_S \left( \left| \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i) - L_D(h) \right| > \epsilon \right) \leq \frac{\delta}{|\mathcal{H}|}$$

assuming,  $0 \leq \ell(h, \mathbf{z}) \leq 1$ . apply Hoeffding's inequality.

\* Hoeffding's inequality.

Let  $\theta_1, \theta_2, \dots, \theta_m$  be a set of iid R.V.s with  $\mu = \mathbb{E}[\theta_i]$ ,

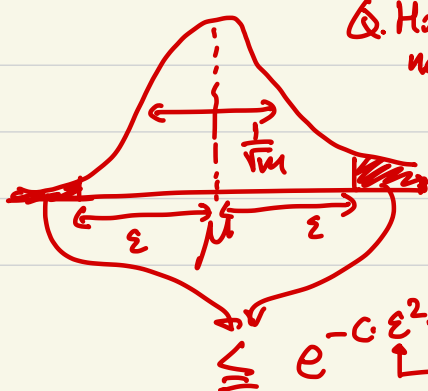
with  $a \leq \theta_i \leq b, \forall i \in [m]$ . then for any  $\epsilon > 0$ ,

$$\mathbb{P} \left( \left| \frac{1}{m} \sum_{i=1}^m \theta_i - \mu \right| > \epsilon \right) \leq 2 \cdot e^{-\frac{2m\epsilon^2}{(b-a)^2}}$$

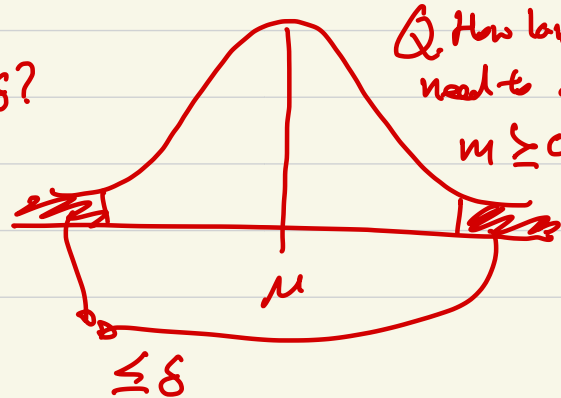
\* first example of concentration of measure.

for fixed  $m$ , varying  $\epsilon$ .

for varying  $m$



Q. How large does  $\epsilon$  need to be for tail  $\leq \delta$ ?  
 $\epsilon = O\left(\frac{1}{\sqrt{m}} \log \frac{1}{\delta}\right)$



Q. How large does  $m$  need to be for tail  $\leq \delta$ ?  
 $m \geq O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta}\right)$

letting

$$\theta_i = l(h, z_i)$$

$$0 \leq l(h, z_i) \leq 1$$

a

b

$$\mu = \mathbb{E}[l(h, z)] = L_0(h)$$

$$\text{Hoeffding's} \Rightarrow \mathbb{P}(|L_S(h) - L_0(h)| > \epsilon) = \mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m l(h, z_i) - L_0(h)\right| > \epsilon\right)$$

$$\leq 2 \cdot e^{-2m\epsilon^2} \stackrel{\text{need}}{\leq} \frac{\delta}{2m}$$

$$\text{need: } e^{-2m\epsilon^2} \leq \frac{\delta}{2m}$$

$$\text{need: } m \geq \left\lceil \frac{\log\left(\frac{2m}{\delta}\right)}{2\epsilon^2} \right\rceil$$

### \* Concentration of Measure.

- The more you know about the R.V.  $\theta_i$ , the tighter concentration you can prove.

① [least information, Markov's Inequality]

if  $\theta_i \geq 0$ ,  $\mathbb{E}[\theta_i] = \mu$  then  $\mathbb{P}(\theta_i \geq \epsilon) \leq \frac{\mu}{\epsilon}$

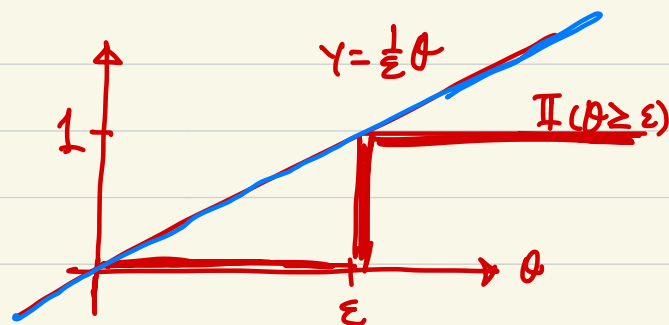
and  $\mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \theta_i \geq \epsilon\right) \leq \frac{\mu}{\epsilon} = O(1)$  : no concentration

proof  $\mathbb{P}(\theta \geq \epsilon) = \int_{\epsilon}^{\theta} f_{\theta}(t) dt$

$$= \mathbb{E}[\mathbb{I}_{\epsilon} \theta \geq \epsilon \}]$$

$$\theta \geq 0 \rightarrow \leq \mathbb{E}\left[\frac{1}{\epsilon} \theta\right]$$

$$\mathbb{E}[\theta] = \mu = \frac{\mu}{\epsilon}$$



② [ 2nd order statistic , Chebyshev's inequality ]

If  $\mathbb{E}[\theta] = \mu$  then  $\mathbb{P}(|\theta_i - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}$   
 $\text{Var}[\theta] \leq \sigma^2$

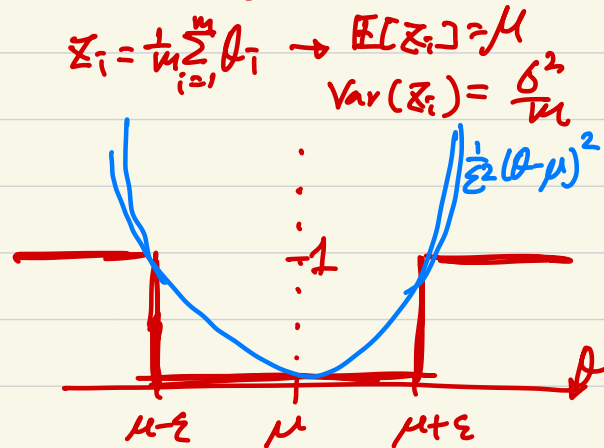
and  $\mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| \geq \varepsilon\right) \leq \frac{\sigma^2}{m\varepsilon^2} = o\left(\frac{1}{m}\right)$   
 ↑ why? Concentration

Proof >

$$\mathbb{P}(|\theta_i - \mu| \geq \varepsilon) = \mathbb{E}[\mathbb{I}_{\{\theta \leq \mu - \varepsilon \text{ or } \theta \geq \mu + \varepsilon\}}]$$

$$\leq \mathbb{E}\left[\frac{1}{\varepsilon^2} (\theta - \mu)^2\right]$$

$$\text{Var}(\theta) \leq \sigma^2 \rightarrow = \frac{\sigma^2}{\varepsilon^2}$$



③ [ bounded domain , Hoeffding's inequality ]

If  $a \leq \theta_i \leq b$  then  $\mathbb{E}\left[e^{\frac{\lambda(\theta_i - \mu)}{m}}\right] \leq e^{\frac{\lambda^2(b-a)^2}{8m^2}}$  (\*)

and  $\mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| \geq \varepsilon\right) \leq 2 \cdot e^{-\frac{2m\varepsilon^2}{(b-a)^2}} = o(e^{-m\varepsilon^2})$

Generic Recipe

⇒ Hoeffding's.

$$\mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \theta_i - \mu > \varepsilon\right) = \mathbb{E}\left[\mathbb{I}_{\left\{\frac{1}{m} \sum_{i=1}^m \theta_i - \mu > \varepsilon\right\}}\right]$$

$$\leq \exp\left\{-\lambda\varepsilon + \frac{\lambda^2(b-a)^2}{8m}\right\}$$

minimize over  $\lambda \geq 0$

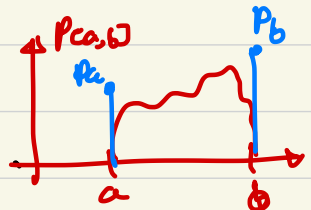
$$\frac{\partial}{\partial \lambda} = -\varepsilon + \frac{\lambda(b-a)^2}{4m} \rightarrow \leq \exp\left\{-\frac{2\varepsilon^2 m}{(b-a)^2}\right\}$$

$$\lambda^* = \frac{4\varepsilon m}{(b-a)^2}$$

\*Repeat for  $\mathbb{P}\left\{\frac{1}{m} \sum_{i=1}^m \theta_i < \mu - \varepsilon\right\}$

We are left to show (\*)  $\mathbb{E}\left[e^{\frac{\lambda(\theta-\mu)}{m}}\right] \leq e^{-\frac{\lambda^2(b-a)^2}{8m^2}}$

$$\mathbb{E}\left[e^{\frac{\lambda(\theta-\mu)}{m}}\right] \leq \max_{P_{[a,b]}} \mathbb{E}_{P_{[a,b]}}\left[e^{\frac{\lambda(\theta-\mu)}{m}}\right]$$



Convexity of  $\exp(\cdot)$   $\leq \max_{P_a, P_b} P_a e^{\frac{\lambda(a-\mu)}{m}} + P_b e^{\frac{\lambda(b-\mu)}{m}}$

s.t.  $P_a + P_b = 1$   
 $aP_a + bP_b = \mu$

$$\leq e^{-\frac{\lambda\mu}{m}} \cdot \left\{ \max_{P_a, P_b} P_a e^{\frac{\lambda a}{m}} + P_b e^{\frac{\lambda b}{m}} \right\} \quad \text{omit proof.}$$

$$\leq e^{-\frac{\lambda\mu}{m}} \cdot \left\{ \max_{P_a, P_b} \exp\left\{\frac{\lambda a}{m} P_a + \frac{\lambda b}{m} P_b + \frac{(b-a)^2 \lambda^2}{8m^2}\right\} \right\}$$

$$= e^{-\frac{\lambda\mu}{m}} \cdot \left\{ \max_{P_a, P_b} \exp\left[\frac{\lambda}{m} \mu + \frac{(b-a)^2 \lambda^2}{8m^2}\right] \right\}$$

$$= \exp\left\{\frac{(b-a)^2 \lambda^2}{8m^2}\right\}$$

④ [Bernoulli dist, Chernoff bound]

$$\theta_i = \begin{cases} 1 & \text{w.p. } p_i \\ 0 & \text{w.p. } 1-p_i \end{cases}$$

$$\mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \theta_i > (1+\delta) \underbrace{\frac{1}{m} \sum_{i=1}^m p_i}_{\mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m \theta_i\right]}\right) \leq e^{-\frac{\delta^2}{2+\frac{2}{3}\delta} \cdot \underbrace{\sum_{i=1}^m p_i}_{O(m)}}$$

$\frac{\delta^2}{2+\frac{2}{3}\delta}$   
 $\downarrow$   
 $O(\delta^2)$   
 or  $O(\delta)$

independent but not

necessarily identically distributed.

$$\mathbb{E}\left[\frac{1}{m} \sum_{i=1}^m \theta_i\right] = \frac{1}{m} \sum_{i=1}^m p_i$$

multiplicative bound

vs.

additive bound  
 $\frac{1}{m} \sum_{i=1}^m p_i + \varepsilon$

Recap.

lecture 10. (ch 2 & 3)

- PAC learning, ← realizable,  $\mathcal{Y} = \{0, 1\}$ ,  $|\mathcal{H}| < \infty$
- technique: union bound.

lecture 11. (ch 4)

- Agnostic PAC learning ← <sup>bounded</sup> non-realizable, any loss,  $|\mathcal{H}| < \infty$
- technique: Hoeffding's ineq.

Q 3.3  $\mathcal{X} = \mathbb{R}^2$ ,  $\mathcal{Y} = \{0, 1\}$ ,  $\mathcal{H} = \{h_r : r \in \mathbb{R}_+\}$ ,  $|\mathcal{H}| = \infty$ ,  
 $h_r = \mathbb{1}_{\{|x| \leq r\}}$

$$m_{\mathcal{H}} = \frac{\log \frac{1}{\delta}}{\epsilon} \rightarrow \text{How?}$$

---

Chapter 5.

the reason the above problem has low sample complexity for PAC learning is that we had a prior knowledge.

- realizable: ground truth  $D, f$  has a circular boundary  
 $\mathcal{H}$  is a collection of circular decision boundaries.

Q. is such prior knowledge necessary for PAC learning? **Yes**.

Learning Problem

fix  $\mathcal{X}, \mathcal{Y}$

nature  $D \xrightarrow{S} \text{Learner } \mathcal{H}, \text{ Algo}$

↓  
 $h$

Goal:  $L_D(h) \downarrow$

Q. Can there be a Universal Algorithm that is PAC learnable for all  $D$ ?

Q. Is prior knowledge necessary for PAC learning?

- lec 10: realizable w.r.t.  $H$ .

- lec 11:  $\min_{h \in H} L_D(h)$  is small (implicit in Agnostic PAC)

- Q3.3: parametric family  $D$

Q. How much prior knowledge is good to assume?

a lot  $\uparrow$  - learner knows  $D$   $\rightarrow h^*$ : Bayes Optimal Predictor

- learner knows a class

- lec 10: realizable,  $|H| < \infty \rightarrow$  PAC  $m \approx \frac{\log |H|}{\epsilon}$

- lec 11: Q3.3,  $f \in H = \{h_r\} \rightarrow$  PAC  $m \approx \log 1/\delta / \epsilon$

non-realizable,  $|H| < \infty \rightarrow$  Agnostic PAC  $m \approx \log \frac{|H|}{\delta} / \epsilon^2$

- learner knows nothing  $\rightarrow$  [No free lunch theorem]

$\downarrow$   
NO