

Relaxing the realizability assumption

Data generating distribution \mathcal{D} is over \mathcal{X}, \mathcal{Y} : $\mathcal{Z}=(x,y) \sim \mathcal{D}$

True error: $L_{\mathcal{D}}(h) \triangleq \mathbb{P}_{\mathcal{D}}(h(x) \neq y)$

Empirical error: $L_S(h) \triangleq \frac{|\{i: i \in [m]: h(x_i) \neq y_i\}|}{m}$

Goal: to find h that minimizes $L_{\mathcal{D}}(h)$

Bayes Optimal Predictor has lower error than any other predictor [HW1 Prob 2]

$$f_{\mathcal{D}}(x) = \begin{cases} 1 & \text{if } \mathbb{P}_{\mathcal{D}}(y=1|x) \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

but requires knowledge of \mathcal{D} .

Agnostic PAC learnability.

A class of hypotheses, \mathcal{H} , is **Agnostic PAC learnable**

if $\exists m_{\mathcal{H}}$ and a learning algorithm s.t. for $\epsilon, \delta, \mathcal{D}$ running Algo on $m \geq m_{\mathcal{H}}(\epsilon, \delta)$ samples returns h s.t. w.p $\geq 1-\delta$

$$L_{\mathcal{D}}(h) \leq \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon$$

Relaxing Binary Classification

Multi-class classification $\mathcal{Y} = [K]$

Regression, $h(x,y) = \mathbb{I}(h(x) \neq y)$, $\mathcal{Y} \subseteq \mathbb{R}$

loss $l_f(x,y) = (h(x) - y)^2$

Risk function $L_{\mathcal{D}}(h) \triangleq \mathbb{E}_{\mathcal{D}} [l_f(h, (x,y))]$

Empirical Risk function $L_S(h) \triangleq \frac{1}{m} \sum_{i=1}^m l_f(h, (x_i, y_i))$

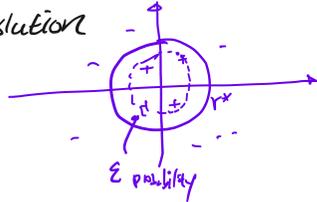
Q3.3 $\mathcal{X} = \mathbb{R}^2, \mathcal{Y} = \{0, 1\}, \mathcal{H} = \{h_r: r \in \mathbb{R}_+\}$

$$h_r(x) = \mathbb{I}(\|x\| \leq r)$$

Prove that \mathcal{H} is PAC learnable (assuming realizability), with

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\ln(1/\delta)}{\epsilon} \right\rceil$$

Solution



fix $\mathcal{D}, f_{\mathcal{D}}$, define $\mathcal{H}_{\text{Bad}} = \{r: r \leq r^*, \mathcal{D}(\|x\| \leq r) \geq \epsilon\}$

Algo: return smallest circle that includes all +.

$$\mathbb{P}_{\mathcal{D}}(\text{Algo} \leq r^*) \leq (1-\epsilon)^m \leq e^{-\epsilon m}$$

for $m \geq \frac{1}{\epsilon} \ln(1/\delta) \leq 6$.

Chapter 4 Uniform Convergence is sufficient for learnability

Definition (ϵ -representative sample)

Training set S is ϵ -representative if

$$\forall h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

Lemma. Assume we are given a training set S that is $\frac{\epsilon}{2}$ -representative. Then, any output of ERM $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$, satisfies

$$L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon$$

Proof strategy: Show ϵ -representative \Rightarrow ERM is ϵ -representative. Use Lemma on S to show ERM is ϵ -representative.

Proof. by triangular inequality $\therefore a \leq b + |a-b|$

$$\begin{aligned} L_D(h_S) &\leq |L_D(h_S) - L_S(h_S)| + L_S(h_S) \\ &\leq \frac{\epsilon}{2} + \underbrace{L_S(h_S) - L_S(h)}_{\leq 0 \sim \text{ERM}} + L_S(h) \\ &\leq \frac{\epsilon}{2} + 0 + |L_S(h) - L_D(h)| + L_D(h) \\ &\leq \frac{\epsilon}{2} + 0 + \frac{\epsilon}{2} + L_D(h) \quad \text{for } \forall h \in \mathcal{H} \end{aligned}$$

To show ERM is agnostic PAC Learner, we are left to show that with probability $1-\delta$, S is ϵ -representative.

Definition (Uniform Convergence)

A hypothesis class \mathcal{H} has the uniform convergence property,

if $\exists m_{\mathcal{H}}^{UC}(\epsilon, \delta)$ s.t. for every ϵ, δ and every D , if $m \geq m_{\mathcal{H}}^{UC}(\epsilon, \delta)$

the S is ϵ -representative with probability at least $1-\delta$.

Corollary 1. \mathcal{H} is agnostically PAC learnable with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)$$

And ERM is the agnostic PAC Learner.

Next, we want to show finite \mathcal{H} is agnostic PAC learnable.
 \uparrow show uniform convergence.

Claim. Any finite \mathcal{H} , i.e., $|\mathcal{H}| < \infty$, is uniformly convergent with

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

Corollary 2. ERM is agnostic PAC learner with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \underbrace{m_{\mathcal{H}}^{UC}\left(\frac{\epsilon}{2}, \delta\right)}_{\text{Corollary 1.}} \leq \underbrace{\left\lceil \frac{\log \frac{2 \cdot |\mathcal{H}|}{\delta}}{\epsilon} \right\rceil}_{\text{claim}}$$

We are left to prove the Claim.

= Q. What is the sample size that guarantees, assuming $\mathbb{E}[\ell(h, \mathcal{Z})] \leq 1$
 $|L_S(h) - L_D(h)| \leq \epsilon$ w.p. $1-\delta$?

n.t.s (need to show) $\mathbb{P}_S(\{\epsilon_S: \forall h \in \mathcal{H}, |L_S(h) - L_D(h)| \leq \epsilon\}) \geq 1-\delta$

n.t.s $\mathbb{P}_S(\{\epsilon_S: \exists h \in \mathcal{H}, |L_S(h) - L_D(h)| > \epsilon\}) < \delta$
 $\leq \sum_{h \in \mathcal{H}} \mathbb{P}_S(\{\epsilon_S: |L_S(h) - L_D(h)| > \epsilon\})$ (*)
Union bound

$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i)$ ← average of iid R.V. $\mathbf{z}_i = (x_i, y_i)$

$L_D(h) = \mathbb{E}_{\mathbf{z} \sim D}[\ell(h, \mathbf{z})] = \mathbb{E}[L_S(h)]$ ← expectation of each term

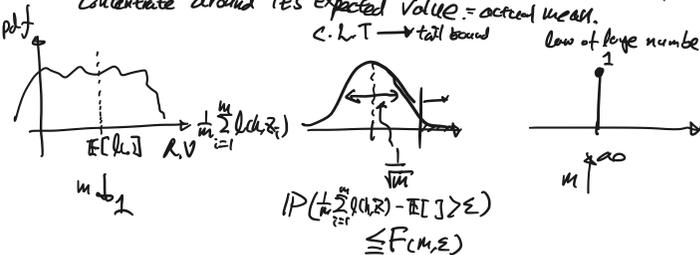
We need to show that the measure of R.V. $L_S(h)$
= distribution

concentrates around its mean.

* most important tool in statistical analysis for learning = concentration of measure
tail bound

Q. How fast (in terms of # of samples) does the measure of empirical mean concentrate around its expected value = actual mean.

C.L.T → tail bound Law of large numbers



* Hoeffding's Inequality

Let $\theta_1, \dots, \theta_m$ be a set of i.i.d random variables with $\mu \equiv \mathbb{E}[\theta_i]$, satisfying

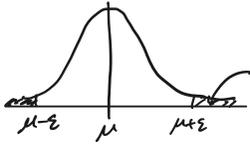
$$\mathbb{P}(a \leq \theta_i \leq b) = 1, \text{ for all } i \in [m] \quad [\text{Boundedness}]$$

then for any $\epsilon > 0$,

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

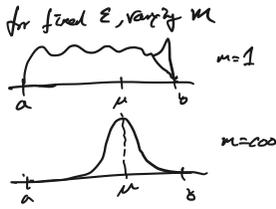
first example of concentration of measure / tail bound

for fixed m , varying ϵ



$$e^{-m\epsilon^2 \cdot C} \quad \text{constant } C = \frac{2}{(b-a)^2}$$

for fixed ϵ , varying m



Letting $\theta_i = f(h, z_i)$, and assuming $a \leq f(h, z_i) \leq b$, and $\mu = \mathbb{E}_D(f(h))$

$$\mathbb{P}_S(|L_S(h) - L_D(h)| > \epsilon) = \mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \epsilon\right)$$

$$\text{Hoeffding's} \rightarrow \leq 2 \cdot e^{-2m\epsilon^2}$$

$$(*) \Rightarrow \mathbb{P}_S(\exists h \in \mathcal{H} : |L_S(h) - L_D(h)| > \epsilon) \leq \sum_{h \in \mathcal{H}} 2e^{-2m\epsilon^2} = 2|\mathcal{H}| \cdot e^{-2m\epsilon^2}$$

$$\text{letting } m(\epsilon, \delta) \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\epsilon^2}, \leq \delta$$

Summary

if $m \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\epsilon^2}$ then \mathcal{H} has uniform convergence

Lemma: Uniform Convergence $\Leftrightarrow m_{\mathcal{H}}^{uc}(\epsilon, \delta) \geq m_{\mathcal{H}}(\epsilon, \delta)$
 (ϵ, δ) -PAC learnable with ERM

$$\text{if } m \geq \left\lceil \frac{2 \cdot \log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{\epsilon^2} \right\rceil \geq m_{\mathcal{H}}^{uc}(\epsilon/2, \delta) \geq m_{\mathcal{H}}(\epsilon, \delta)$$

then \mathcal{H} is (ϵ, δ) -PAC learnable with ERM.

End of Lecture 3.

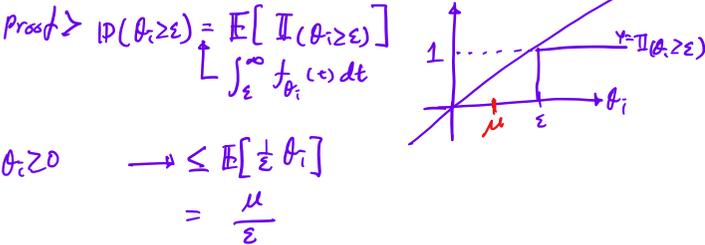
Concentration of measure

The more you know about the R.V. θ_i , the tighter concentration you can prove.

① [least information, Markov's inequality]

$$\text{If } \theta_i \geq 0 \text{ w.p.1 then } \mathbb{P}(\theta_i \geq \epsilon) \leq \frac{\mu}{\epsilon} \quad \text{no concentration, why?}$$

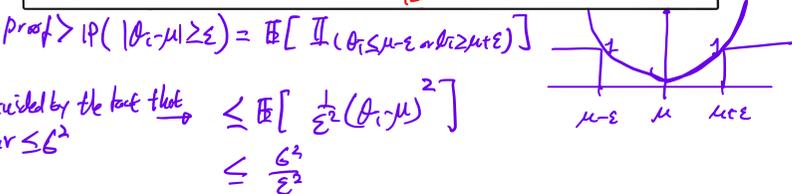
$$\text{and } \mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \theta_i \geq \epsilon\right) \leq \frac{\mu}{\epsilon} = O\left(\frac{1}{\sqrt{m}}\right) \quad \text{right?}$$



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

$$\text{If } \text{Var}(\theta_i) \leq \sigma^2 \text{ then } \mathbb{P}(|\theta_i - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$$

$$\text{and } \mathbb{P}\left(\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| \geq \epsilon\right) \leq \frac{\sigma^2}{m\epsilon^2} = O\left(\frac{1}{\sqrt{m}}\right) \quad \text{variance } \times \frac{1}{m}$$



③ [bounded domain, Hoeffding's inequality]

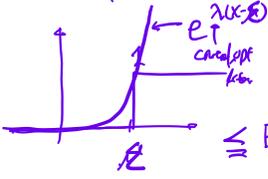
$$\text{If } a \leq \theta_i \leq b \text{ then } \mathbb{E}\left[e^{\lambda(\theta_i - \mu)}\right] \leq e^{\frac{\lambda^2(b-a)^2}{8m^2}} \quad (**)$$

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

Generic Recipe

Special to Hoeffding's

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



$$\leq \mathbb{E}\left[e^{\lambda\left(\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right) - \varepsilon}\right]$$

$$\leq e^{-\lambda\varepsilon} \prod_{i=1}^m e^{\lambda\left(\theta_i - \mu\right)}$$

$$\leq e^{-\lambda\varepsilon} \prod_{i=1}^m e^{\frac{\lambda^2 (b-a)^2}{8m^2}}$$

opt over $\lambda \rightarrow$

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

$$-\varepsilon + \frac{2\lambda(b-a)^2}{8m} = 0 \rightarrow \leq \exp\left(-\frac{2m\varepsilon^2}{(b-a)^2}\right)$$

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

do the same for $\frac{1}{m} \sum_{i=1}^m \theta_i < \mu - \varepsilon$

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

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

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

Part 1: $a=P_a, b=P_b$

$$\leq e^{-\frac{\lambda\mu}{m}} \max_{P_a, P_b} P_a e^{\frac{\lambda a}{m}} + P_b e^{\frac{\lambda b}{m}}$$

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

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

* Chernoff's Bound

$\theta_1 \dots \theta_m$ independent Bernoulli dist. $\theta_i \in \{0, 1\}$ w.p. P_i and $1-P_i$
but not identically distributed

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

multiplicative

2 phase

$$\geq \frac{\delta^2}{2+\frac{2}{3}\delta}$$