

7. Learning via Uniform Convergence

Tuesday, January 18, 2022 11:28 AM

§4.1[SS] Uniform Convergence

A sufficient condition for PAC learning

(necessary in some cases, cf. Thm. 6.7, though modern perspective is that uniform convergence is often asking for too much)

Def A training set S is an ε -representative sample (with respect to set $Z = X \times Y$, set \mathcal{H} , loss ℓ , distr. D) if

$$(\forall h \in \mathcal{H}) \quad |\hat{L}_S(h) - L_D(h)| \leq \varepsilon$$

hence "uniform" flavor, e.g. $\|\hat{L}_S - L_D\|_\infty \leq \varepsilon$

Such a set is quite nice:

Lemma 4.2 If S is $\varepsilon/2$ -representative then if $h_S \in \text{ERM}_{\mathcal{H}}(S)$ then

$$L_D(h_S) = \left(\min_{h \in \mathcal{H}} L_D(h) \right) + \varepsilon \quad \text{i.e. } \in \arg\min_{h \in \mathcal{H}} \hat{L}_S(h)$$

exactly what we want for PAC learning

Proof:

let $h \in \arg\min_{h \in \mathcal{H}} L_D(h)$

$$\begin{aligned} L_D(h_S) &\leq \hat{L}_S(h_S) + \varepsilon/2 && \leq \hat{L}_S(h) + \varepsilon/2 \quad [\text{since } h_S \text{ minimizes } \hat{L}_S] \\ &\stackrel{\text{by } \varepsilon/2\text{-repr.}}{\leq} L_D(h) + \varepsilon/2 + \varepsilon/2 \quad [\text{by } \varepsilon/2\text{-repr. again}] \\ &= \left(\min_{h \in \mathcal{H}} L_D(h) \right) + \varepsilon \end{aligned} \quad \square$$

Building on this concept...

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

(w.r.t. to Z, ℓ) if $\exists m_{\mathcal{H}}^{\text{UC}} : (0, 1)^2 \rightarrow \mathbb{N}$ s.t.

$\forall \varepsilon, \delta \in (0, 1)$, \forall prob. distr. D over $Z = X \times Y$,

if S has $m \geq m_{\mathcal{H}}^{\text{UC}}(\varepsilon, \delta)$ iid samples, then

S is ε -representative with probability $\geq 1 - \delta$

Uniform over
 $h \in \mathcal{H}$ and
all possible
 D

Corollary If \mathcal{H} has the uniform convergence property with $m_{\mathcal{H}}^{\text{UC}}$,
then ① \mathcal{H} is agnostic PAC learnable

$$\textcircled{2} \quad m_H(\varepsilon, S) \leq m_H^{\text{uc}}(\varepsilon, S)$$

\textcircled{3} ERM_H is a valid algorithm to learn H
(proof is immediate)

§4.2 Finite classes ($|H| < \infty$) have the uniform conv. property

Fix ε, δ . We want S to be ε -representative w.p. $\geq 1 - \delta$,

i.e. $D^m(\{S : \forall h \in H, |\hat{L}_S(h) - L_D(h)| \leq \varepsilon\}) \geq 1 - \delta$

or $D^m(\{S : \exists h \in H \text{ s.t. } \dots > \varepsilon\}) \leq \delta$

$\underbrace{\dots}_{\text{event } E}$

$$E = \bigcup_{h \in H} \{S : |\hat{L}_S(h) - L_D(h)| > \varepsilon\}$$

Same trick
we did last
lecture

so via union bound, $D^m(E) \leq \sum_{h \in H} D^m(E_h)$

Rmk: union bound
aka Boole's Ineq.

So look at $D^m(E_h)$

(Before, in realizable case, $\hat{L}_S(h) = 0$ and we could
break $D^m(E_h)$ into $\prod_{i=1}^m D(\dots)$ via independence...
this trick no longer applies)

Since $\hat{L}_S(h) = \frac{1}{m} \sum_i l(h, z_i)$ and $z_i \sim D$ iid,

we have $\mathbb{E} \underbrace{\hat{L}_S(h)}_{\text{r.v.}} = \underbrace{L_D(h)}_{\text{its mean}}$

So today's lecture is statistics: how likely is a r.v. to deviate
significantly from its mean?

= Statistics Interlude =

1) as $m \rightarrow \infty$, the law of large numbers says

$$\hat{L}_{S(m)}(h) \rightarrow L_D(h), \text{ but there's no rate}$$

2) Central Limit Thm says $\sqrt{m} \hat{L}_S(h) \sim \text{Normal}$ as $m \rightarrow \infty$

(basic CLT has no rate, but Berry-Esseen CLT does)
Could we use facts about normal distr.?

No, not good bounds: see Vershynin book

3) What about Chebyshev's Inequality?

Recall Markov's Ineq: if $Z \geq 0$ is a r.v. then

$$\mathbb{P}[Z \geq a] \leq \frac{\mathbb{E}[Z]}{a}$$

Let θ be a r.v.
w/ variance σ^2 , then Chebyshev's Ineq. says
 $\forall \varepsilon > 0, \quad \mathbb{P}[|\theta - \mathbb{E}[\theta]| > \varepsilon] \leq \frac{\sigma^2}{\varepsilon^2}$

let's see if we can apply to our situation

Define $\theta_i = l(h, z_i)$, a r.v. w/ mean $L_D(h)$

$$\text{so } \theta = \hat{L}_S(h) = \frac{1}{m} \sum \theta_i$$

What is $\text{Var}(\theta)$? Suppose $\text{Var}(\theta_i) = \sigma^2$

$$\begin{aligned} \text{Var}\left(\frac{1}{m} \sum \theta_i\right) &= \frac{1}{m^2} \text{Var}\left(\sum \theta_i\right) \\ &= \frac{1}{m^2} \cdot \sum \text{Var}(\theta_i) \quad \text{by independence} \\ &= \frac{1}{m} \cdot \sigma^2 \end{aligned}$$

So via Chebyshev, $\mathbb{P}[|\hat{L}_S(h) - L_D(h)| > \varepsilon] \leq \frac{\sigma^2}{m \varepsilon^2}$

That's for event E_h , so now union bound

$$D^m(E) \leq \sum_{h \in H} D^m(E_h) \leq |H| \cdot \frac{\sigma^2}{m \varepsilon^2} =: \delta$$

So our sample complexity m_H^{uc} would be

$$m_H^{uc}(\varepsilon, \delta) = \frac{1}{\varepsilon^2} \cdot \frac{1}{\delta} \cdot \sigma^2 \cdot |H|$$

not great but workable not good vague but acceptable
terrible:
we want/need a $\log(|H|)$ or similar

Ex: $H = \{ \text{1D model, parameter is double precision number} \}$

$$\text{so } |H| = 2^{64}$$

then $|H|$ vs $\log(|H|)$ is a BIG DEAL

Union bound + Chebyshev = Bad Idea

4) Chernoff-style bounds: Hoeffding, Chernoff, Bernstein etc
 "Concentration Ineq." Exploit \sum iid structure

Many variants - see Handout.

We'll do Hoeffding

First, some tools

Def A Rademacher (or symmetric Bernoulli) r.v. X
 means $X = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases}$ usually $\{0, 1\}$ output

and a Rademacher vector means each component is an independent Rademacher

Lemma "Exponentiated Markov"

If X is a r.v., and $\lambda \in \mathbb{R}$ a constant, $\xrightarrow{\lambda > 0}$ aka Moment Generating Functn
 $\mathbb{P}[X \geq \varepsilon] \leq \exp(-\lambda\varepsilon) \cdot \mathbb{E}[e^{\lambda X}]$ MGF(λ)

proof

$$\begin{aligned}\mathbb{P}[X \geq \varepsilon] &= \mathbb{P}[e^{\lambda X} \geq e^{\lambda\varepsilon}] \quad \text{since } X \mapsto e^{\lambda X} \text{ is} \\ &\leq e^{-\lambda\varepsilon} \mathbb{E}[e^{\lambda X}] \quad \text{via Markov's} \quad \text{monotonically increasing } (\lambda > 0) \\ &\quad (e^{\lambda\varepsilon} \gg 0)\end{aligned}\quad \square$$

Lemma "Hoeffding's Lemma"

If X is Rademacher, $\mathbb{E}[e^{\lambda X}] \leq e^{\lambda^2/2}$

proof

$$\mathbb{E}[e^{\lambda X}] = \frac{1}{2}e^\lambda + \frac{1}{2}e^{-\lambda} = \cosh(\lambda) \leq e^{\lambda^2/2} \quad \text{via Taylor Remmehbr Thm. or similar.} \quad \square$$

now...

Thm (Hoeffding, simple, one-sided)

Let X_i be i.i.d Rademacher, $\vec{a} \in \mathbb{R}^n$, then $\forall \varepsilon > 0$,

$$\mathbb{P}\left[\sum_{i=1}^n a_i X_i \geq \varepsilon\right] \leq \exp\left(-\frac{\varepsilon^2}{2\|\vec{a}\|_2^2}\right) \quad \text{Common cases:} \\ a_i = 1 \Rightarrow \|\vec{a}\|_2^2 = n \\ a_i = \frac{1}{n} \Rightarrow \|\vec{a}\|_2^2 = \frac{1}{n}$$

Thm (Hoeffding, more general)

Let X_i be independent and bounded $X_i \in [m_i, M_i]$, then $\forall \varepsilon > 0$,

$$\mathbb{P}\left[\left|\sum_{i=1}^n (X_i - \mathbb{E}[X_i])\right| \geq \varepsilon\right] \leq 2 \cdot \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (M_i - m_i)^2}\right)$$

proof:

Let $\lambda > 0$ (TBD). wlog let $\|\vec{a}\|_2 = 1$

By exponentiated Markov,

$$\begin{aligned}\mathbb{P}\left[\sum_i a_i X_i \geq \varepsilon\right] &\leq \exp(-\lambda\varepsilon) \mathbb{E}\left[\exp(\lambda \sum_i a_i X_i)\right] \\ &= \exp(-\lambda\varepsilon) \mathbb{E}\left[\prod_i \exp(\lambda a_i X_i)\right] \quad \text{how exp works} \\ &= \exp(-\lambda\varepsilon) \prod_i \mathbb{E}[\exp(\lambda a_i X_i)] \quad \text{independence} \\ &\leq \exp(-\lambda\varepsilon) \prod_i \exp\left(\frac{(\lambda a_i)^2}{2}\right) \quad \text{Hoeffding lemma} \\ &= \exp(-\lambda\varepsilon) \exp\left(\lambda^2/2 \sum_i a_i^2\right) \quad \text{how exp works} \\ &= \exp(-\lambda\varepsilon + \lambda^2/2) \quad \text{since } \|\vec{a}\|_2^2 = 1\end{aligned}$$

now choose λ (to minimize the bound) polynomial

i.e., $\lambda = \varepsilon$, so ...

$$\dots = \exp(-\varepsilon^2/2) \quad \square$$

- Ex $X_i \stackrel{iid}{\sim} \text{Rademacher}$, $i \in [m]$, so $E[X_i] = 0$
 $E[X_i^2] = 1$ hence $\sigma^2 = 1$
- 1) Hoeffding $P\left[\left|\frac{1}{m} \sum_{i=1}^m X_i\right| \geq \varepsilon\right] \leq 2 \exp(-m\varepsilon^2/2)$ i.e. $\vec{a} \in \mathbb{R}^m$, $a_i = \frac{1}{m}$, $\|\vec{a}\|_2^2 = \frac{1}{m}$
 - 2) Chebyshov $P\left[\left|\frac{1}{m} \sum_{i=1}^m X_i\right| \geq \varepsilon\right] \leq \frac{\sigma^2}{m} \cdot \frac{1}{\varepsilon^2} = \frac{1}{m\varepsilon^2}$ since $\text{Var}\left(\frac{1}{m} \sum_i X_i\right) = \frac{\sigma^2}{m}$
 - 3) Exact: by symmetry, $P\left[\left|\frac{1}{m} \sum_i X_i\right| \geq \varepsilon\right] = 2 \cdot P\left[\frac{1}{m} \sum_i X_i \geq \varepsilon\right]$
 transform: $X_i = \begin{cases} 1 & \Rightarrow \tilde{X}_i = \begin{cases} 1 & \\ 0 & \end{cases} \text{ so } \tilde{X}_i = \frac{1}{2}(X_i + 1) \\ -1 & \end{cases}$
 $\tilde{X}_i = 2\tilde{X}_i - 1$
 $\text{so } = 2 \cdot P\left[\sum_{i=1}^m (2\tilde{X}_i - 1) \geq m \cdot \varepsilon\right]$
 $= 2 \cdot P\left[\sum_i \tilde{X}_i \geq \frac{m}{2}(1 + \varepsilon)\right]$
 $= 2 \cdot (1 - P\left[\sum_i \tilde{X}_i < \frac{m}{2}(1 + \varepsilon)\right])$
 $= 2(1 - F(\frac{m}{2}(1 + \varepsilon)))$ where $F(x)$ is CDF of
 $\alpha \sim (\frac{1}{2}, m)$ binomial r.v.

<u>Ex</u>	<u>True</u>	<u>Hoeffding</u>	<u>Chebyshov</u>
$m = 2000, \varepsilon = 0.1$	$6.8 \cdot 10^{-6}$	$9.1 \cdot 10^{-5}$.05
$m = 100,000, \varepsilon = 0.02$	$2.5 \cdot 10^{-10}$	$4.1 \cdot 10^{-9}$.025

... going back to our uniform convergence analysis ...

Hoeffding
 So via Chebyshov, $P\left[\left|\hat{L}_s(h) - L_0(h)\right| \geq \varepsilon\right] \leq \frac{2 \exp(-m\varepsilon^2/2)}{m\varepsilon^2}$
 That's for event E_h , so now union bound

$$D^m(E) \leq \sum_{h \in H} D^m(E_h) \leq |H| \cdot 2 \cdot \exp(-m\varepsilon^2/2) =: \delta$$

i.e., solving for m , $-m\varepsilon^2/2 = \log(\delta/|H|)$

$$m = \frac{2/|H|}{\varepsilon^2} \log\left(\frac{2/|H|}{\delta}\right)$$

Summary:

Corollary 4.6 U.C. for finite classes (binary classification, agnostic)

Let $|H| < \infty$, l be the 0-1 loss function, then h

has the uniform convergence property with

$$m_H^{uc}(\varepsilon, \delta) \leq \lceil \frac{1}{2\varepsilon^2} \log\left(\frac{2|H|}{\delta}\right) \rceil$$

... and hence agnostic PAC learnable with

$$m_H(\varepsilon, \delta) = m_H^{uc}(\varepsilon/2, \delta)$$

my derivation above
 differs by a factor
 of 2, but not a
 big deal
 (don't trust me...
 but don't trust
 book either)