

Rates of Convergence in ERM

We've proved learnability for opt hypo spaces

Next, we'll derive quantitative results on rate of convergence for $\mathcal{L}_{\text{sample}}$. Which is crucial to prove uniform convergence / PAC-learnability & control capacity of hypo space.

(1) For finite hypo space:

Lemma. For \mathcal{H} finite, $\hat{\mu}_n$ \mathcal{H} -ERM learner for unbiased ERF $\hat{\mathcal{L}}_n$. If $Z_{v,n} := \mathcal{L}(\mu)(v) - (\mathcal{L}_n(\hat{\mathcal{L}}_n(v)) + \mathcal{L}_n(\mu))$ is subgaussian for $\forall v \in \mathcal{H}$ with var. σ_n^2 (indep of \mathcal{H}).

Then: i) $\mathbb{E}(\sup_n Z_{v,n}) \leq (2\sigma_n^2 \log |\mathcal{H}|)^{\frac{1}{2}}$.

ii) if $Z_{v,1}$ is sub-gaussian with var.

$$\sigma_1^2 = \sigma, \quad \hat{\mathcal{L}}_n = \frac{1}{n} \sum_{j=1}^n \mathcal{L}(x_j | v), \quad \mathcal{L}_n = \frac{1}{n}$$

and $\mathcal{L}_n(\mu) = \mathcal{L}(\mu)$ in addition \Rightarrow

$$\sigma_n^2 = \sigma^2/n, \quad \mathbb{E}(\sup_n Z_{v,n}) \leq (2\sigma^2 \log |\mathcal{H}|/n)^{\frac{1}{2}}$$

iii) If $\mu \in \mathcal{M}$. $X_j \overset{i.i.d.}{\sim} \mu$. $\exists V_n \in \mathcal{M}$. $Z_{V,n} = 0$
 with ii) $\Rightarrow \mathbb{E}(\mathcal{L}(\mu \| \hat{\mu}_n)) \leq 2 \left(\frac{2\sigma^2 \log |\mathcal{M}|}{n} \right)^{\frac{1}{2}}$

Cor. $\mathbb{P}(\mathcal{L}(\mu \| \hat{\mu}_n) \geq \epsilon) \leq 2 \left(\frac{2\sigma^2 \log |\mathcal{M}|}{n \epsilon^2} \right)^{\frac{1}{2}}$

$$\mathbb{P}(\mu \neq \hat{\mu}_n) \leq 2 \left(2\sigma^2 \log |\mathcal{M}| / n \right)^{\frac{1}{2}} / \inf_{\mathcal{M}/\{\mu\}} \mathcal{L}(\mu \| \nu)$$

Pf: By Markov inequality. and $|\mathcal{M}|$ finite

$$\mathbb{P}(\mu \neq \hat{\mu}_n) = \mathbb{P}(\mathcal{L}(\mu \| \hat{\mu}_n) \geq \inf_{\mathcal{M}/\{\mu\}} \mathcal{L}(\mu \| \nu))$$

Pf: i) LHS = $\mathbb{E}(\alpha^{-1} \log e^{\alpha \sup_{\mathcal{M}} Z_{v,n}})$

$$\stackrel{\text{Jensen}}{\leq} \alpha^{-1} \log \mathbb{E}(e^{\alpha \sup_{\mathcal{M}} Z_{v,n}})$$

$$\leq \alpha^{-1} \log \mathbb{E}\left(\sum_{\mathcal{M}} e^{\alpha Z_{v,n}}\right)$$

$$\stackrel{\text{subgauss}}{\leq} \alpha^{-1} \log \mathbb{E}\left(\sum_{\mathcal{M}} e^{\frac{1}{2} \sigma_n^2 \alpha^2}\right)$$

$$= \alpha^{-1} \log(|\mathcal{M}| e^{\frac{1}{2} \sigma_n^2 \alpha^2}) = \frac{\sigma_n^2 \alpha}{2} + \frac{\log |\mathcal{M}|}{\alpha}$$

choose $\alpha = (2 \log |\mathcal{M}|)^{\frac{1}{2}} / \sigma_n$ to optimize.

ii) Note only $\mathcal{L}(X_k | \mathcal{V})$ is random.

iii) From RMK of \mathcal{L} we decompose:

$$\mathbb{E}(\mathcal{L}(\mu \| \hat{\mu}_n)) = 2 \mathbb{E}_{\mathcal{M}}(\sup_{\mathcal{M}} Z_{v,n}) + \mathbb{E}_{\text{mod}}(\mu).$$

Def: Given metric \mathcal{L} on \mathcal{M} . $\bar{\mathcal{L}} > 0$. Empirical process $\{Z_v\}_{v \in \mathcal{M}}$ is subgaussian process with

Var. $\bar{L}^2 \bar{\lambda}^{-2}$ if $\mathbb{E}(Z_V) = 0$. $\forall V \in \mathcal{K}$ and $Z_n - Z_V$ is subgaussian with var. $\bar{L}^2 \bar{\lambda} (p, V)^2$.

Remark: i) It gives Z_n kind of Lip. conti.

which is crucial when considering c.p.b. hypo space \subset in balls

ii) $\bar{\lambda}$ isn't necessary to spri. div. d

which measures success of learning

\subset some can't be metric. eg. KL. even

eg. Note under cond. of Lem ii) above. We

have $Z_{V,n} = \frac{1}{n} \sum_{i=1}^n \ell(X_i | V) - \mathbb{E}_V \ell(X | V)$

We claim: $Z_V := Z_{V,1}$ is subgaussian for

$V \in \mathcal{K} = \{N(m, 1) \mid m \in [-1, 1]\}$. $\bar{\lambda} \in N(m, 1), N(m', 1)$

$= |m - m'|$. $X \sim N(m, 1)$. where $\ell(X | V)$

$= \log f_V(x) = \frac{1}{2} (x - m)^2 + \log \sqrt{2\pi}$, $V = N(m, 1)$.

Since $\mathbb{E}_X \ell(X | N(m, 1)) = \frac{1}{2} (m - m)^2 + \frac{1}{2} + \log \sqrt{2\pi}$

$\Rightarrow Z_n = \frac{1}{2} (X - m)^2 - \frac{1}{2} (m - m_0)^2 - \frac{1}{2}$

$\therefore Z_n - Z_{m'} = 2(m' - m)X + \text{const.}$ which

is Gaussian r.v.

Next, we want to control number of center of covering balls.

Def: For $\varepsilon > 0$, \bar{d} metric on \mathcal{X} . $(N, D(\varepsilon)) \stackrel{\Delta}{=} N, D(\varepsilon, \bar{d}, \mathcal{X})$

i) Covering number $N(\varepsilon)$ is smallest number of center v_j . s.t. $\mathcal{X} \subset \bigcup_{j=1}^{N(\varepsilon)} \bar{B}(v_j, \varepsilon)$

ii) Packing number $D(\varepsilon)$ is largest number of $v_j \in \mathcal{X}$. s.t. $\min_{i,j} \bar{d}(v_i, v_j) > \varepsilon$.

Lemma: $D(2\varepsilon) \stackrel{i)}{\leq} N(\varepsilon) \stackrel{ii)}{\leq} D(\varepsilon)$.

Pf: i) For $v_j, j=1, \dots, N(\varepsilon)$. $\bigcup_{j=1}^{N(\varepsilon)} \bar{B}(v_j, \varepsilon) \supset \mathcal{X}$.

And $\mu_j, j=1, \dots, D(2\varepsilon)$, packing pts.

At most one μ_j can be contained in some ε -ball $\bar{B}(v_i, \varepsilon)$ since its diameter is at most 2ε . So: $D(2\varepsilon) \leq N(\varepsilon)$

ii) $v_j \in \mathcal{X}, j=1, \dots, D(\varepsilon)$. If $\{\bar{B}(v_j, \varepsilon)\}$ can't cover \mathcal{X} . i.e. $\exists v \notin \bigcup \bar{B}(v_j, \varepsilon)$.

We can set $v_{D(\varepsilon)+1} = v \Rightarrow$ contradiction.

Def: For $(\mu) \in \mathbb{R}^k$ cpt with $(\mu) \in B_r(0)$. We

have $\mathcal{X} = \{\mu_k\}_{0 \leq k \leq \infty}$ and $\bar{d} = d_{\mu}$ satisfies:

$$N(\varepsilon) \leq D(\varepsilon) \leq \left(1 + \frac{2r}{\varepsilon}\right)^A.$$

Pf: For $V_j, j=1, \dots, D(\varepsilon)$, $\bar{B}(V_j, \frac{\varepsilon}{2})$ disjoint

balls. s.t. $\bigcup \bar{B}(V_j, \frac{\varepsilon}{2}) \subset B_{r+\frac{\varepsilon}{2}}(\theta')$

$$\Rightarrow D(\varepsilon) \cdot C_A \left(\frac{\varepsilon}{2}\right)^A \leq C_A \left(r + \frac{\varepsilon}{2}\right)^A$$

$$\text{i.e. } D(\varepsilon) \leq \left(1 + \frac{2r}{\varepsilon}\right)^A.$$

Th. (Dudley's inequality on max entropy estimate)

$\{z_v\}_\mathcal{N}$ is subgaussian with var. proxy $\bar{L}\bar{\alpha}^2$

$\mathcal{H}(\mathcal{N}, \bar{L})$ is b.h., separable metric space.

$v \mapsto z_v(w)$ is conti. for IP-a.s. w . Then:

$$\mathbb{E} \left(\sup_{\mathcal{N}} z_v \right) \leq 12 \sum_{k=0}^{\infty} 2^{-k} \left(\log(N(2^{-k}, \bar{L}\bar{\alpha}, \mathcal{N})) \right)^{\frac{1}{2}}$$

Where $k_0 \in \mathbb{Z}^+$ is largest s.t. $N(2^{-k_0}, \bar{L}\bar{\alpha}, \mathcal{N}) = 1$.

$$\underline{\text{Cor.}} \quad \mathbb{E} \left(\sup_{\mathcal{N}} z_v \right) \leq 24 \int_0^\infty \left(\log(N(\varepsilon, \bar{L}\bar{\alpha}, \mathcal{N})) \right)^{\frac{1}{2}} d\varepsilon$$

$$= 24 \bar{L} \int_0^\infty \left(\log(N(\varepsilon, \bar{\alpha}, \mathcal{N})) \right)^{\frac{1}{2}} d\varepsilon$$

$$\underline{\text{Pf:}} \quad 12 \sum_{k=0}^{\infty} \square = 24 \sum_k \int_{2^{-(k+1)}}^{2^{-k}} \left(\log(N(2^{-k}, \dots)) \right)^{\frac{1}{2}}$$

$$\leq 24 \sum_{k=0}^{\infty} \int_{2^{-(k+1)}}^{2^{-k}} \left(\log(N(\varepsilon, \bar{L}\bar{\alpha}, \mathcal{N})) \right)^{\frac{1}{2}} d\varepsilon.$$

$$\leq 24 \int_0^\infty \left(\log(N(\varepsilon, \bar{L}\bar{\alpha}, \mathcal{N})) \right)^{\frac{1}{2}} d\varepsilon.$$

And $N(\varepsilon, \bar{L}\bar{\alpha}, \mathcal{N}) = N(\frac{\varepsilon}{\bar{L}}, \bar{\alpha}, \mathcal{N})$. let

$\varepsilon = \varepsilon/\bar{L}$. substitution.

Remark: For ε large enough. $N(\varepsilon, \bar{\lambda}, \mathcal{K}) = 1$ So:

$\int_N^\infty \square = 0$ will not explode while $\int_0^1 \square$ part is dangerous.

Pf: i) If $N(\varepsilon, \bar{\lambda}, \mathcal{K}) = \infty$ for ε small. Then:

It's nothing to prove. WLOG. $N(\varepsilon, \bar{\lambda}, \mathcal{K}) < \infty$

By conti. $\sup_{\mathcal{K}} Z_V = \sup_{\mathcal{K}'} Z_V$ a.s. where \mathcal{K}'

is countable dense set of \mathcal{K} .

WLOG. let $\mathcal{K} = \{h_n\}_{n \in \mathbb{Z}^+}$ countable set.

ii) Set $\mathcal{K}_s := \{h_n\}_{n \in \mathbb{N}_s}$. Choose $V_{k,j}$ of \mathcal{K}_s

$\bar{\lambda}$ -balls with $r = 2^{-k}$ covering \mathcal{K}_s .

let k_1 satisfy $\forall k \geq k_1, N(2^{-k}, \bar{\lambda}, \mathcal{K}_s) = s$

then: $\{V_{k,j}\} = \{h_n\}_{n \in \mathbb{N}_s}$ in this case.

And k_0 satisfies $\forall k \leq k_1, N(2^{-k}, \bar{\lambda}, \mathcal{K}_s) = 1$.

For $V \in \mathcal{K}_s$. let $Z_k(V) = V_{k,j}$ where

$V \in \bar{B}^{\bar{\lambda}}(V_{k,j}, 2^{-k}) \Rightarrow Z_{k_1}(V) = V, Z_{k_0}(V) = V_{k_0,1}$

So: $Z_V = \sum_{k_0}^{k_1} (Z_{Z_k(V)} - Z_{Z_{k-1}(V)}) + Z_{Z_{k_0}(V)} = Z_{V_{k_0,1}}$

$\Rightarrow E(\sup_{\mathcal{K}_s} Z_V) \leq \mathbb{I}_{k_0}^{k_1} E(\sup_{\mathcal{K}_s} (Z_{Z_k(V)} - Z_{Z_{k-1}(V)}))$

3) Over \mathcal{K}_s , they're at most M many

different $z_{2^k(v)} - z_{2^{k+1}(v)}$, where

$$m = N(2^{-k}, \bar{L}, \bar{\lambda}, \mathcal{K}_s) \cdot N(2^{-(k+1)}, \bar{L}, \bar{\lambda}, \mathcal{K}_s)$$

$$\leq D(\dots, \mathcal{K}_s) D(\dots, \mathcal{K}_s) \stackrel{\text{mono.}}{\leq} D(\dots, \mathcal{K}) D(\dots, \mathcal{K})$$

$$\stackrel{\text{Lem.}}{\leq} N(2^{-(k+1)}, \bar{L}, \bar{\lambda}, \mathcal{K})^2.$$

With $\bar{L}, \bar{\lambda}(z_{2^k(v)}, z_{2^{k+1}(v)})$

$$\leq \bar{L}(\bar{\lambda}(z_{2^k(v)}, v) + \bar{\lambda}(z_{2^{k+1}(v)}, v)) \stackrel{\text{def}}{\leq} 3 \cdot 2^{-k}$$

Apply Lem.: $\mathbb{E}(\sup_{\mathcal{K}_s} z) \leq$

$$\sum_{k=1}^{K_1} \left(2 \cdot 3 \cdot 2^{-k} \right)^2 \log(N(2^{-(k+1)}, \bar{L}, \bar{\lambda}, \mathcal{K}))^2)^{\frac{1}{2}}$$

Then use MCT. let $s \uparrow \infty$.

cor. Under condition of Thm above where $z_{v,n}$

$$= \lambda(\mu|v) - (c_n \bar{L}_n(v) + h_n(\mu)) \text{ with } \bar{L}_n(v)$$

$$= \frac{1}{n} \sum_{j=1}^n \ell(X_j|v) \text{ unbiased, } c_n = \frac{1}{n}, X_k \stackrel{\text{i.i.d.}}{\sim} \mu, z_v \stackrel{\text{def}}{=} z_{v,1}$$

Then: for ERM-learner $\hat{\mu}_n$, if $\exists v^* \in \mathcal{K}$ s.t.

$$z_{v^*,n} = 0. \Rightarrow \text{we have:}$$

$$\mathbb{E}(\lambda(\mu|\hat{\mu}_n)) \leq 48 \bar{L} n^{\frac{1}{2}} \int_0^\infty (\log(N(\varepsilon, \bar{L}, \mathcal{K}))^2)^{\frac{1}{2}} d\varepsilon + \Sigma_{n, \text{mod}}$$

And for $\mu \in \mathcal{J} \subset \mathcal{K} \Rightarrow \Sigma_{n, \text{mod}} = 0$.

$$S_\varepsilon = \mathbb{P}(\lambda(\mu|\hat{\mu}_n) > \varepsilon) \leq \frac{48 \bar{L}}{\sqrt{n} \varepsilon} \int_0^\infty (\log(N(\varepsilon, \bar{L}, \mathcal{K}))^2)^{\frac{1}{2}} d\varepsilon.$$

Rmk: In prior example, we have: $N(\epsilon)$

$\leq (1 + 2^{r/2})^L \Rightarrow RNS$ is finite!

Besides, RNS doesn't depend on n .

So for $\mathcal{J} = \mathcal{H}$, it's PAC-learnable
w.r.t. L_{KL} and by \hat{I}_n here.

Pf: From \mathcal{Z}_V to $\mathcal{Z}_{n,V}$. We only need to
replace L by L/\sqrt{n} . Apply them above.