

Empirical Risk Min.

$\mathcal{L} : \mathcal{H} \ni v \mapsto \mathcal{L}(p||v)$ is loss func. and \mathcal{H} is set of all p.m. candidates for learning

Def: i) \mathcal{H} is called hypothesis space.

Rmk: Choice of \mathcal{H} is related to data
If we want to learn the list. p
by choosing v from \mathcal{H} . The best choice
is $v \in \arg \min_{v \in \mathcal{H}} \mathcal{L}(v)$. But often we
don't know real list. p .

So we need $\tilde{\mathcal{L}}_n(\cdot)$ empirical risk func.
Computed from data $X = (X_1, \dots, X_n)$. And
it should be good approx. for $\mathcal{L}(\cdot)$.

Rmk: i) $\mathcal{L}(v)$ measure the degree of failure
from choice $v \in \mathcal{H}$.

ii) For \mathcal{H} opt. $v \in \mathcal{H} \mapsto \mathcal{L}(p||v)$ is
cont. $\Rightarrow \arg \min_{\mathcal{H}} \mathcal{L}(v) \neq \emptyset$.

For \mathcal{H} is liv. $\Rightarrow \arg \min_{\mathcal{H}} \mathcal{L}(v)$

is singleton: $d(\mu||v) \leq d(\tilde{\mu}||v) +$

$d(\mu||\tilde{\mu}) \Rightarrow d(\mu||\tilde{\mu}) = 0 \Rightarrow \mu = \tilde{\mu}$

ii) Empirical risk func. $\hat{L} := \{\hat{L}_n\}$ for

given \mathcal{X} . \mathcal{F} is family of map \hat{L}_n :

$\mathcal{Y}_n \times \mathcal{X}^n \rightarrow \mathbb{R}$. st.

a) $\mathcal{X}^n \ni x_n \mapsto \hat{L}_n(v, x_n)$ is measurable

for $\forall v \in \mathcal{Y}_n$ $\arg \min \hat{L}_n(v, x_n)$

b) $\exists c_n > 0$ and $h_n: \mathcal{Y} \rightarrow \mathbb{R}$. s.t. $\exists C_n$

for $\forall v \in \mathcal{Y}_n$, $\forall k$. $\forall n \in \mathcal{Y}$. we have

$c_n \hat{L}_{nk}(v, x_{nk}) + h_n(v) \xrightarrow{pr} L(\mu||v)$

provided $x_k \stackrel{i.i.d.}{\sim} \mu$.

rank: If $\mathcal{H}_n \subset \mathcal{T}$. Then: we can replace (c_n)

by (1) and consider $h \rightarrow 0$.

iii) ERF is unbiased if $\forall \mu \in \mathcal{Y}$. $\forall v \in \mathcal{Y}_n$

$x_j \stackrel{i.i.d.}{\sim} \mu$. $E_{\mu}[\hat{L}_n(v, x_n) + h_n(v)] = L(\mu||v)$

iv) SLA $\hat{\mu}_n(x_n)$ is called empirical risk

minimizer (ERM)-learner if:

$\hat{\mu}_n(x_n) \in \arg \min_{\mu_n} \hat{L}_n(v, x_n)$.

and the process of minimizing $\hat{L}(\cdot, x_n)$
 is called training of ERM-learner.

Next we omit n and focus on $\mathcal{H} = \mathcal{H}_n$.
 We want to refine ERM-learner w.r.t. the
 i.i.d. data model (similar for DTMC)

(1) Max. Likelihood estimate as ERM:

Def: $V \in \mathcal{M}_1^+(\mathbb{R}^d)$, p.m. $X_n = (X_1, \dots, X_n)$ is n-sample
 of i.i.d. r.v. X_k with values in \mathbb{R}^d .

i) $\mathcal{H} \subset \mathcal{M}_1^+(\mathbb{R}^d)$, hypothesis space s.t. its
 elements have discrete or conti. density.

$$\text{ii)} \hat{L}_n(x_n | V) = \begin{cases} \hat{\pi}_V(V \in \{x_i\}), & V \text{ is discrete} \\ \hat{\pi}_V f_V(x_i), & V(dx) = f_V(x)dx. \end{cases}$$

for $V \in \mathcal{H}$.

Recall MLE is SLA sc. it fulfills:

$$\hat{\mu}_n(x_n) \in \arg \max \{ \hat{L}_n(x_n | V) \mid V \in \mathcal{H} \}.$$

Remark: We want to minimize $-\log \hat{L}_n(x_n | V)$
 Next $\log f_V(x_i) \& \log V(x_i) = \ell(x_i | V)$.

Lem. $X_j \stackrel{i.i.d.}{\sim} \mu \in \mathcal{I} \subset M^+(X^1)$. X_n is n -observation

Assume $\tilde{\pi} \ll \mathcal{I} \subseteq \text{p.m.}$ with conti. density

$$\text{i.e. } V(x) = f_V(x) \lambda_x \text{ & } L(x, v) = \log f_V(x) \in L'(\mu)$$

for $\theta, v \in \tilde{\pi}$ and $\mu \in \mathcal{I}$.

$$\text{Then: } \hat{L}_n(v, X_n) = -\log \tilde{L}(X_n, v) = -\sum \hat{L}(x_j | v)$$

is unbiased ERF w.r.t. d_{KL} with $c_n = \frac{1}{n}$

$$\text{and } h_n(\mu) = \mathbb{E}_\mu \hat{L}(x | \mu).$$

Rmk: It's similar to prove for discrete case.

Pf: Unbiased is from i.i.d. para $[x_j]$.

$$\text{And by SLLN: } c_n \hat{L}_n(v, X_n) =$$

$$-\frac{1}{n} \sum \hat{L}(x_j | v) \rightarrow -\mathbb{E}_\mu \hat{L}(x | v)$$

$$\Rightarrow c_n \hat{L}_n(v, X_n) + h_n(\mu) \rightarrow -\mathbb{E}_\mu \log \frac{\tilde{\pi}_v}{\pi_\mu}$$

(2) Error Decomp.:

Thm. For i.i.d model $X_k \sim \mu$ and λ divergence

$L(\mu || \lambda) = h(\mu || \lambda)$. \mathcal{H}_n is hypo space and

$\{\hat{L}_n\}$ is unbiased ERF w.r.t. L with c_n, h_n

If μ_n is a SLA. Then we have:

$$i) 0 \leq \lambda(\mu || \hat{\mu}_n) \leq \Sigma_{n \text{ mod } \ell} \epsilon_\mu + \Sigma_{n \text{ learn}} + 2 \Sigma_{n \text{ sample}} \epsilon_\mu$$

where $\Sigma_{n \text{ mod}} \epsilon_\mu = \inf_{\mathcal{H}_n} \lambda(\mu || v)$.

$$\Sigma_{n \text{ learn}} = C_n \left(\hat{I}_n(\hat{\mu}_n, X_n) - \inf_{\mathcal{H}_n} \hat{I}_n(v, X_n) \right)$$

$$\Sigma_{n \text{ sample}} \epsilon_\mu = \sup_{\mathcal{H}_n} | \lambda(\mu || v) - (C_n \hat{I}_n(v, X_n) + h_n \epsilon_\mu) |.$$

ii) For ERM-learner $\hat{\mu}_n$, we have:

$$\lambda(\mu || \hat{\mu}_n) \leq \Sigma_{n \text{ mod}} \epsilon_\mu + 2 \Sigma_{n \text{ sample}} \epsilon_\mu.$$

iii) In addition with ii). For $\mu \in \mathcal{T} \subset \mathcal{H}_n$.
we have: $\lambda(\mu || \hat{\mu}_n) \leq 2 \Sigma_{n \text{ sample}} \epsilon_\mu$.

(\mathcal{T} is set of "true" measures to be learned)

Pf: For $v \in \mathcal{H}_n$. Note that

$$\begin{aligned} \lambda(\mu || \hat{\mu}_n) &\stackrel{(*)}{=} \lambda(\mu || v) + [C_n \hat{I}_n(\hat{\mu}_n) + h_n \epsilon_\mu] \\ &\quad - [C_n \hat{I}_n(v) + h_n \epsilon_\mu] \Big] + \{ [\lambda(\mu || \hat{\mu}_n) - (C_n \hat{I}_n(\hat{\mu}_n) + h_n \epsilon_\mu)] \\ &\quad - [\lambda(\mu || v) - (C_n \hat{I}_n(v) + h_n \epsilon_\mu)] \Big] + [C_n \hat{I}_n(\hat{\mu}_n) - \hat{I}_n(v)] \Big\}. \end{aligned}$$

$$\leq \lambda(\mu || v) + C_n (\hat{I}_n(\hat{\mu}_n) - \hat{I}_n(v)) + 2.$$

$$\sup_{v \in \mathcal{H}_n} | \lambda(\mu || v) - (C_n \hat{I}_n(v) + h_n \epsilon_\mu) |$$

Take $\inf \{ \dots | v \in \mathcal{H}_n \}$ on RHS.

Rmk: i) If $\exists V_n^* \in \mathcal{H}_n$. s.t. $L(\mu || V_n^*) = \inf_{V_n} L(\mu || V_n)$

\Rightarrow We can modify above:

$$0 \leq L(\mu || \hat{\mu}_n)$$

$$\leq L(\mu || U) + C_n (\hat{L}_n(\hat{\mu}_n) - \hat{L}_n(V))$$

$$+ [C_n \hat{L}_n(V) + h_n(\mu)] - L(\mu || U)$$

$$+ [L(\mu || \hat{\mu}_n) - C_n \hat{L}_n(\hat{\mu}_n) - h_n(\mu)]$$

$$\leq [C_n + \sup_{V_n} |L(\mu || U) - (C_n \hat{L}_n(V) + h_n(\mu))|]$$

And we take $U = V_n^*$ or RHS. \Rightarrow

We have a better estimate. Note that

One term $\xrightarrow{pr} 0$ by def of ERFs.

We only need to control $\sup \square$

i) For $\{\hat{L}_n\}$ unbiased and $\hat{\mu}_n$ ERM.

in case i). We take $\bar{E}(\cdot)$ on the

Result i): Since $\hat{L}_n(\hat{\mu}_n) \leq \hat{L}_n(V_n^*)$

$$\therefore \bar{E}(L(\mu || \hat{\mu}_n)) \stackrel{\text{unbiased}}{=} L(\mu || V_n^*) +$$

$$\bar{E}(\sup_{V \in \mathcal{H}_n} |L(\mu || U) - (C_n \hat{L}_n(V) + h_n(\mu))|)$$

(By Chebyshev Ineq. we can estimate

$$P(L(\mu || \hat{\mu}_n) > \varepsilon) \text{ or } P(L(\mu || \hat{\mu}_n) - L(\mu || V_n^*) > \varepsilon)$$

\Rightarrow If $N \in \mathcal{F}$. To get learnability of n :

a) $\mu \in \bar{\mathcal{N}} := \overline{V_n \mathcal{N}_n}^k$

b) $\sup_{\mu_n} |E_{\mu_n}(\hat{\mu}_n) - (E_n \hat{L}_n(\hat{\mu}_n) + h_n(\mu_n))| \rightarrow 0$

(Note for pointwise v. b) holds by the def. of EMF $\{\hat{L}_n\}$)
if $\mu \in \mathcal{F} \subset \mathcal{N}$. then $E_{\mu}(\hat{\mu}_n) \rightarrow 0$. and
we only need cond. b).

i) Note $\hat{\mu}_n = \hat{\mu}_{n \in \mathcal{X}_n}$ repeat on \mathcal{X}_n . So:

$E(E_n \hat{L}_n(\hat{\mu}_n) + h_n(\mu_n)) \neq E_{\mu}(\hat{\mu}_n)$ in
general. \Rightarrow Take expectation on (*) in
proof may not give better estimate.

ii) When capacity $N \in \mathcal{F} \Rightarrow \sum_{n \in \mathcal{N}} c_{\mu_n} \downarrow$ but
 $\sum_{n \in \text{sample}} c_{\mu_n}$ (and $E_{\mu_n} c_{\mu_n}$) \uparrow .

v) $E_{\mu_n}(\hat{\mu}_n) - (E_n \hat{L}_n(\hat{\mu}_n) + h_n(\mu_n))$ can be
large when $\hat{L}_n(\hat{\mu}_n) \downarrow$ by increasing
the data. It's known as overfitting
 \Rightarrow To decouple it. we take suprema so
let $\hat{\mu}_n$ disappear.

Cor. Under cond. of Thm iii) above with

\mathcal{H} finite. $\rightarrow \forall \mathcal{S} \subseteq \mathcal{H}$ is PAC-learnable.

$$\begin{aligned} \text{pf: } & \text{If } \mathbb{P}(\mathcal{L}_\mu || \hat{\mu}_n) > L \xrightarrow{\text{Learn}} \mathbb{P}(\mathcal{L}_{\sup_{\mathcal{H}}} | \mathcal{D}) > \varepsilon \\ & = \mathbb{P}\left(\bigcup_{\mathcal{H}} \left[\mathcal{L}_\mu || \hat{\mu}_n - \left(C_n \tilde{L}_n + \ln(\mu) \right) \right] > \frac{\varepsilon}{2} \right), \\ & \leq \sum_{\mathcal{H}} \mathbb{P}(|\dots| > \frac{\varepsilon}{2}) \rightarrow 0 \end{aligned}$$

So we can define $\kappa(\varepsilon, \delta) = \text{max}$

$$\left[\mu_{\mu, v} < \frac{\varepsilon}{2}, \frac{\delta}{L_{\mathcal{H}}} \right] \mid m \in \mathcal{S}, v \in \mathcal{V} \}.$$

Rmk: Define μ is agnostically learnable

if $\mathcal{L}_\mu || \hat{\mu}_n - \text{Err}_{\text{ind}}(\mu) \rightarrow 0$ in pr.

(When $\text{Err}_{\text{ind}}(\mu) \rightarrow 0$ then we put

it on LHS. And PAC case is similar)

So if $\mathcal{S} \neq \mathcal{H}$. \mathcal{S} is still agnostically PAC-learnable.

(3) Opt hypo space:

$|\mathcal{H}| < \infty$ is too strong and restrictive in 1st Cor. So we want to consider infinite hypo. space.

Q1. Consider parametric space $\Theta = [\mu_-, \mu_+]$
 $\times [\sigma_-, \sigma_+]$. $M_\theta = N(\mu, \sigma)$ for $\theta = (\mu, \sigma)$

$L(\mu_\theta, \mu_*) = \| \theta - \theta^* \|_2$. It's infinitely uncount.

Next, we assume ERM-learner ($\Sigma_{n, \text{learn}} = 0$)
and $\mathcal{T} \subseteq \mathcal{H}$ ($\Sigma_{n, \text{model}} = 0$). Thus, we only need

$$\Sigma_{n, \text{sample}} = \sup_n | \hat{L}_n(v) + h(\mu) - L(\mu)(v) | \rightarrow 0 \text{ in pr.}$$

i.i.d.
for $\mu \in \mathcal{T}$ and $x_k \sim \mu$.

Consider unbiased ERM has form: $(\hat{L}_n(v)) =$
 $\frac{1}{n} \sum_i^n \ell(x_i | v)$ if $v \in \mathcal{H}$ is uni. distributed

$$d\mu(x) = f(x|v)dx$$

Note that there're two questions:

1) Whether suprema is measurable. i.e. Is $\Sigma_{n, \text{sample}}$
a measurable r.v.?

It's true if \mathcal{H} is λ -separable. then $\exists \mathcal{H}$.

dense countable & $v \mapsto \ell(x|v)$ conti. \Rightarrow

$$\Sigma_{n, \text{sample}} = \sup_{\mathcal{H}_n} |\cdot|.$$

2) Show $\Sigma_{n, \text{sample}} \xrightarrow{\text{pr}} 0$

We require some uniform LLN. It can

be achieved by assuming \mathcal{H} is opt.

\Rightarrow we can consider $|c_n \bar{\lambda}^n(\mu) + h(\mu) - \lambda \|\mu\|_V| \xrightarrow{pr} 0$

in every small ball by conti.: $V \mapsto L(x|V)$

Thm. (Uniform L-LN)

$\exists \subset \mathcal{H} \subset M^+(\mathbb{R}^d)$. \mathcal{H} is λ -opt. For $x_k \sim \mu \in \mathcal{J}$. i.i.d. data mod. If $V \mapsto L(x|V)$ is conti. and $k(x) = \sup_{\mu} |L(x|\mu)| \in L^\infty$. Then:

Example $c_{\mu^n} \rightarrow 0$. a.s. (\Rightarrow the ERM learner

$\hat{\mu}_n$ learns $\mu \in \mathcal{J}$)

Pf: We first prove:

$$\lim_{n \rightarrow \infty} \overline{\sup} \frac{1}{n} \sum_j L(x_j|V) \leq \sup_n \overline{E}(L(x|V)) = 1$$

Set $\varphi(x|V, \epsilon) = \sup_{V' \in \mathcal{V}, \|V-V'\|_V \leq \epsilon} L(x|V') \leq k(x).$

$\varphi(x|V, \epsilon) \downarrow L(x|V)$ by conti. of L .

With DCT: $\overline{E}(\varphi(x|V, \epsilon)) \downarrow \overline{E}(L(x|V))$

\Rightarrow for $\Sigma > 0$. $\exists \epsilon_V$. s.t.

$$\overline{E}(\varphi(x|V, \epsilon_V)) \leq \overline{E}(L(x|V)) + \Sigma.$$

$\Rightarrow \exists \{V_i\}_{i \in I}$ covers \mathcal{H} .

By def., $\exists K$. St. $\frac{1}{n} \sum_{j=1}^n \ell(x_j | v) \leq \frac{1}{n} \sum_{j=1}^n \ell(x_j | v_k - e_{v_k})$

take \sup_K and then $\sup_{v \in \mathcal{V}}$:

$$\frac{1}{n} \sum_{j=1}^n \ell(x_j | v_k - e_{v_k}) \leq \frac{1}{n} \sum_{j=1}^n E[\ell(x_j | v_k)] + \epsilon$$

$$\rightarrow E[\ell(x | v_k)] + \epsilon. P\text{-a.s.}$$

$$S_0 := \lim_{n \rightarrow \infty} \sup_{1 \leq k \leq r} \frac{1}{n} \sum_{j=1}^n \ell(x_j | v_k - e_{v_k}) \leq \sup_{k \in \mathcal{V}} E[\ell(x | v_k)] + \epsilon.$$

Let $\epsilon \rightarrow 0$. We have:

$$\overline{\lim}_{n \rightarrow \infty} \sup_{v \in \mathcal{V}} \frac{1}{n} \sum_{j=1}^n \ell(x_j | v) \leq \sup_{v \in \mathcal{V}} E[\ell(x | v)]. P\text{-a.s.}$$

Next, we see $\ell(x | v)$ is conti. So that

$\overline{E}[\ell(x | v)]$ is conti. Set $\bar{\ell}(x | v) = \overline{\ell}(x | v)$

$- \overline{E}[\ell(x | v)]$. Repeat above on $\{\bar{\ell}(x | v)\}$.

$$\Rightarrow \overline{\lim}_{n \rightarrow \infty} \sup_{v \in \mathcal{V}} \frac{1}{n} \sum_{j=1}^n \bar{\ell}(x_j | v) \leq 0. P\text{-a.s.}$$

Apply on $-\bar{\ell}(x | v)$. We have $\underline{\lim}_{n \rightarrow \infty} \inf \bar{\ell}(x | v) \leq 0$.

Besides, we note that $\frac{1}{n} \sum_{j=1}^n \bar{\ell}(x_j | v) =$

$$= \frac{1}{n} \sum_{j=1}^n (\ell(x_j | v) + h_n(\mu) - h_n(\hat{\mu}) - \overline{E}[\ell(x | \bar{\ell}(x | v))])$$

$$= n \bar{\ell}(x | v) + h_n - \ell(\mu | v).$$

Remark: i) It's not PAC-learnable because

the regularity condition.

ii) μ can be $\widehat{\tau}$ -opt. $\widehat{\tau}$ is other metric. And assume $V \mapsto L(X|V)$ is $\widehat{\tau}$ -conti. The THm still holds.
(So we can choose weaker topo.
to let μ be opt.)

(4) Consistent Para. MLE:

For model of parametric statistics $\Sigma_{\mu_0} \otimes \Theta$

recall that $\mu = \bar{J}_n(\mu_n(\cdot))$ is $\lambda\Theta$ -opt (\Leftrightarrow)
 $\Theta \subset \mathbb{X}^n$ is 1-1-opt. (if $\mu_n(\cdot)$ is injective)

Rank: μ will also be λ -opt if Θ is 1-1-opt
and $\Theta \ni \theta \mapsto \mu_\theta$ is λ -conti.

Set $\widehat{\theta}_n = \arg \min_{\theta \in \Theta} \widehat{L}_n(\mu_\theta)$. We say consistency

of para. estimate $\{\widehat{\theta}_n\}$ for θ_0 if $|\widehat{\theta}_n - \theta_0| \xrightarrow{\text{pr}} 0$

Rank: We have $\mu_{\widehat{\theta}_n}$ is ERM w.r.t. $\widehat{L}_n(\cdot)$.

Thm. If para. model $\Sigma_{\mu_0} \otimes \Theta = \Sigma_{\mu^*} \otimes \mathbb{X}^n$, s.t. ④
 $\subseteq \mathbb{X}^n$ opt. $\Theta \ni \theta \mapsto \mu_\theta$ is injective &
 λ -conti. Then:

$\hat{M}_n = \mu \hat{\theta}_n$ (errors $J = n$ w.r.t. $\lambda \Leftrightarrow$

$\hat{\theta}_n$ is consistent for Θ .

Cor. So if cond. of uniform LLN holds.

then - the maximal likelihood para.

estimator $\hat{\theta}_n$ is consistent

Pf: By Lem. before, we have :

$$\lambda(\hat{\mu}_n, \mu_0) \sim \lambda(\Theta^*(\hat{\mu}_n, \mu_0)) = (\hat{\theta}_n - \theta_0).$$

Rem: It works for normal / exponential / Poi.

/ binomial / geo. dist. But uniform dist. is special since it isn't even conti. v.r.t KLL.

L'g. Chose $K = KLL$. $\Theta = [\mu_-, \mu_+] \times [\sigma_-, \sigma_+]$.

$$\theta = (\mu, \sigma) \in \Theta \mapsto M_\theta = e^{-\frac{1}{2} \sum (x - \mu)^2 / \sigma^2} / \sqrt{2\pi \sigma^2}$$

is injective. And its KLL-conti. is also easy to check. Also Θ is go.