# Stat. Learning

(1) Basic limit THms:

The i.i.d model is good because of the existence of SLLN.

Next, we consider DTMC model:

Let $S = \{c_1 \cdots c_d\}$. finite set

**Lem.** $\ell = (\ell_s)$ is initial dist. $p = (p_{ss'})$ is transfer matrix. $\Rightarrow$ dist. of $X_t$ is

$$\ell_s^{(t)} := \mathbb{P}(X_t = s) = (\ell^\top p^t)_s.$$

**Pf:** $\ell_s^{(t+1)} \overset{induct}{=} \sum_{s'} p_{s's} \ell_{s'}^{(t)} = (\ell^\top p^t)_s.$

**Def:** i) discrete density $\bar{e}$ on $S$ is called stationary state if $\bar{e}^\top p = \bar{e}^\top$.

ii) A DTMC is strongly mixing if $\exists$ stationary state $\bar{e}$. st. for any initial dist. $\ell$. : $\lim\limits_{t \to \infty} \ell^\top p^t = \bar{e}^\top$.

Rmk: Stationary state for the strongly mixing DTMC is unique.

Next, we let $g : S \to \mathbb{R}^1$. And we want to estimate $\mathbb{E}_{\bar{e}}(g) := \sum_S g(s) \bar{e}(s)$.

Thm. $(X_t)$ is strongly mixing MC on $S$ with initial dist $\ell$. Then:

$$\frac{1}{T} \sum_0^{T-1} g(X_t) \xrightarrow{pr} \mathbb{E}_{\bar{e}}(g) \quad \text{where } \bar{e} \text{ is its}$$

stationary state.

Pf: $\mathbb{P}( | \frac{1}{T} \sum_0^{T-1} g(X_t) - \mathbb{E}_{\bar{e}}(g) | \geq \mathcal{E} ) \leq$

$$\frac{1}{\mathcal{E}^2} \frac{1}{T^2} \sum_{t, \tau'}^{T-1} \mathbb{E}( (g(X_t) - \mathbb{E}_{\bar{e}}(g))(g(X_{t'}) - \mathbb{E}_{\bar{e}}(g)) )$$

$$\overset{A}{=} \frac{1}{T^2 \mathcal{E}^2} \sum_{t, t' \geq 0}^{T-1} A_{t, t'}.$$

$$A_{t, t'} = \sum_{s, s'} [\ell^T p^t]_s (g(s) - \mathbb{E}_{\bar{e}}(g)) P_{s, s'}^{t'-t}$$

$$\cdot (g(s') - \mathbb{E}_{\bar{e}}(g))$$

$$= \sum \square (P_{s, s'}^{t'-t} - \bar{e}_{s'}) \square +$$

$$\underset{=0}{\sum \square \bar{e}_{s'} \square} \quad \text{for } t \leq t'.$$

$$\text{So}: \ |A_{t,t'}| \leq 4\|g\|_a^2 |S| \max_{s,s'} |P_{s,s'}^{t'-t} - \bar{e}_{s'}| \xrightarrow{t'-t\to\infty} 0$$

$$\Rightarrow RHS = \left( \sum_{|t'-t|\geq \log T}^{T-1} + \sum_{|t'-t|\leq \log T}^{T-1} \square \right) / T^2 \varepsilon^2$$

$$\leq C \sum_{|t'-t|\geq \log T}^{T-1} A_{t,t'} + T\log T \Big) / T^2 \varepsilon^2$$

$$\leq \left( T^2 \cdot o(1) + T\log T \right) / T^2 \varepsilon^2 \xrightarrow{T\to\infty} 0$$

$\underline{Cor.}$ It can converge $P$-a.s. if $A_{t,t'}$ satisfy

$$|A_{t,t'}| \leq C^{|t-t'|}. \quad \text{for} \quad c < 1.$$

$\underline{Pf}$: Set $Y_T = \frac{1}{T} \sum_0^{T-1} g(X_t) - E_{\bar{e}}(g)$ and

$$S_T = T Y_T.$$

$$\text{So}: \ \mathbb{P}\big( |Y_{T^2}| \geq \varepsilon \big) \lesssim T^{-2} \to 0 \quad \text{which}$$

follows from prop. of $A_{t,t'}$.

For $m \in [T^2, (T+1)^2)$. We have:

$$\mathbb{P}\big( |S_m - S_{T^2}| \geq T^2 \varepsilon \big)$$

$$\lesssim T^{-4} \Big( \sum_{T^2+1}^{(T+1)^2-1} m - T^2 \Big)$$

$$\text{So}: \ \sum_T \mathbb{P}\big( |S_m - S_{T^2}| \geq T^2 \varepsilon \big) \leq \sum \frac{1}{T^2} < \infty$$

$\underline{Thm.}$ If $P = (P_{s,s'})_{s\times s}$ satisfies $P_{s,s'} > 0. \ \forall s,s'$

Then: the DTMC has a unique stationary

state $\bar{e}$ and satisfies:

$$\sup_{s,s'} |P^t_{s,s'} - \bar{e}_{s'}| \xrightarrow{t\to\infty} 0 \quad \text{exponentially.}$$

Pf: $V := \{ e = (e_s)_s \mid e_s \in [0,1], \sum_s e_s = 1 \}$.

with metric $d(e, e') = \sum_s |e_s - e'_s|$

is complete ( CLS of $(\mathbb{R}^{|S|}, \|\cdot\|_1)$ )

Set $T : V \to V$. $e \mapsto T(e) = (e^T P)^T$

$\varepsilon := \inf \{ P_{s,s'} \mid s,s' \in S \} > 0$.

Next, we apply Banach fixed pt.

$d(T(e), T(e')) = \sum_{s'} | \sum_s (e_s - e'_s) P_{s,s'} |$

$\quad = \sum_{s'} | \sum_s (e_s - e'_s)(P_{s,s'} - \varepsilon) |$

$\quad \leq \sum_{s'} \sum_s | e_s - e'_s | (P_{s,s'} - \varepsilon)$

$\quad = (1 - |S| \varepsilon) \, d(e, e').$

If $\varepsilon = 1/|S|$. Set $\bar{e} = (\frac{1}{|S|}, \cdots, \frac{1}{|S|})$.

(2) Stat. Learning Algorithm:

Next, we consider set $S \subseteq \mathbb{R}^\lambda$. $|S| < \infty$.

Pf: 1) $d(\cdot, \|\cdot\|)$ is divergence on $M^+_1(\mathbb{R}^\lambda)$. the

space of p.m's on $\mathbb{R}^\lambda$. if $d : M^+_1(\mathbb{R}^\lambda)$

$* \mathcal{M}_1^+(\mathcal{X}^\wedge) \to \bar{\mathbb{R}}_{\geq 0}$, s.t. $d(\mu\|\nu) = 0 \iff \mu = \nu$

if $d$ is true metric. We write:

$d(\cdot\|\cdot) = d(\cdot,\cdot)$.

<u>Rmk</u>: $d$ takes role of measuring how

close our estimate $\hat{\mu}_n$ to $\mu$.

ii) Unsupervised stat. learning algorithm is

collection of func. $\{\hat{\mu}_n\}_n$ defined by

$\hat{\mu}_n: \mathcal{X}^{\wedge \times n} \to \mathcal{M}_1^+(\mathcal{X}^\wedge)$. s.t. $\forall \mu \in \mathcal{M}_1^+(\mathcal{X}^\wedge)$.

$\mathcal{X}^{\wedge \times n} \ni (x_1, \cdots x_n) =: \mathcal{X} \longmapsto d(\mu\|\hat{\mu}_n(\mathcal{X}))$ is

measurable from $(\mathcal{X}^{\wedge n}, \mathcal{B}_{\mathcal{X}^\wedge n})$ to $(\bar{\mathbb{R}}_{\geq 0}, \mathcal{B}_{\bar{\mathbb{R}}_{\geq 0}})$.

<u>Rmk</u>: Supervised stat. learning is a input-

output model. But there's no input

$(i.e. \ labeled \ training \ data)$ here.

iii) $X_j: (\Omega, \mathcal{A}, \mathbb{P}) \to (\mathcal{X}^\wedge, \mathcal{B}_{\mathcal{X}^\wedge})$ is data model

$\mathcal{X}_n = (X_1, \cdots X_n)$ is sample of size $n$.

Let $\hat{\mu}_n = \hat{\mu}_n \circ \mathcal{X}_n$ random SLA. $X_j \overset{i.i.d}{\sim} \mu$.

We call $\mu$ is $d$-learnable for the

divergence $d$ if $d(\mu\|\hat{\mu}_n) \overset{pr}{\to} 0$.

<u>Rmk</u>: It's kind of weak learnable.

$\mathcal{J} \subseteq \mathcal{M}_1^+(\varphi^d)$ is called $\alpha$-learnable if $\forall \mu \in \mathcal{J}$ is $\alpha$-learnable.

$\mathcal{J}$ is called $\alpha$-PAC-learnable (probably approxi. correct) if $\exists n(\cdot, \cdot): (0,1) \times (0,1) \to \mathbb{N}$. St. $\forall \mu \in \mathcal{J}$. $X_j \sim \mu$. $\forall \varepsilon < 1, \delta > 0$ $\exists n(\varepsilon, \delta)$. We have:

$$p(\alpha(\mu \| \hat{\mu}_n) > \varepsilon) \leq \delta. \quad \forall n \geq n(\varepsilon, \delta).$$

<u>Rmk</u>: i) PAC-learnable is stronger than $\alpha$-learnable. Since $n(\varepsilon, \delta)$ won't depend on $\mu \in \mathcal{J}$.

ii) Restriction of data generating dist. on $\mathcal{J}$ represent the prior knowledge of data.

iii) Note that $\mathcal{J} \subseteq \overline{\varlimsup_{n} \operatorname{Im}(\hat{\mu}_n)}^d$. But if $\mathcal{J} \neq F$, Set $\operatorname{Im}(\mu)^{\overset{A}{=}}$

$\inf_{\sim} \inf_{\nu \in \operatorname{Im}(\hat{\mu}_n)} \alpha(\mu \| \nu)$ then. $\exists \mu \in \mathcal{J}:$

st. $\sum_{x \in \mathcal{X}} \langle \mu \rangle > 0$. We can replace

$d \langle \mu \| \hat{\mu}_n \rangle \to 0$ by $d \langle \mu \| \hat{\mu}_n \rangle \to \sum_{x \in \mathcal{X}} \langle \mu \rangle$

## ① Learning on DTMC:

Denote: i) $d_s \langle \mu, \nu \rangle = \max_s | \mu(\{s\}) - \nu(\{s\}) |$.

$$\mathbb{I}_{\{ \mu(s) = \nu(s) = 1 \}} + \mathbb{I}_{\{ \mu \neq \nu \} \cap \{ \mu(s) = \nu(s) = 1 \}^c}$$

ii) Empirical measure $\hat{\mu}_{T-1} = \frac{1}{T} \sum_0^{T-1} \delta_{x_t}$

Lemma. $\{ \hat{\mu}_T \}$ learns $\mathcal{J} = \{ \mu \in \mathcal{M}_1^+ \langle \mathcal{X}^d \rangle \mid \mu(s) = 1 \}$

w.r.t. $d_s$ for $\{ X_t \}$ of a strongly

mixing DTMC with invar. measure $\mu$

Pf: Set $\bar{\ell}_s = \mu(\{s\})$. $\hat{\ell}_{s,T} = \hat{\mu}_T(\{s\})$

$$p \langle d_s \langle \bar{\ell}, \hat{\ell}_T \rangle > \varepsilon \rangle \leq$$

$$\sum_s p \langle | \bar{\ell}_s - \hat{\ell}_{s,T} | > \varepsilon \rangle \xrightarrow[\text{mixing}]{T \to \infty} 0 \quad \langle |S| < \infty \rangle$$

## ② Learning on i.i.d model:

Pf: $\mathcal{G} \leq \{ g : \mathcal{X}^d \to \mathcal{X}' \mid g \text{ is bdd. measurable} \}$

$\mathcal{G}$ - weak topo is generated by seminorms

$d_g \langle \mu, \nu \rangle = | \int_{\mathcal{X}^d} g \, d\mu - g \, d\nu |. \quad g \in \mathcal{G}.$

i.e. $\mu_n \xrightarrow{\mathcal{Z}-w} \mu$ if $d_g(\mu, \mu_n) \to 0$. $\forall g \in \mathcal{Z}$.

$\mathcal{Z}$ - strong top. is generated by $\mathcal{Z}$-divergence

$$d_{\mathcal{Z}}(\mu \| \mu_n) := \sup_{\mathcal{Z}} d_g(\mu, \mu_n) \to 0.$$

Rmk: i) We say $\mathcal{Z}$ is separating if:

$d_{\mathcal{Z}}$ is a norm. i.e. $d_{\mathcal{Z}}(\mu, \nu) = 0$

$(\Longrightarrow) \mu = \nu$

ii) $M^+(\mathcal{X})$ isn't LS. But it can

be embedded into sign measure

space which's LS. So the "norm"

will make sense.

iii) The Lem. in ① satisfies Df by

choosing $\mathcal{Z} = [\mathbb{1}_{(-\infty, s]} | s \in S]$.

Next, we want to investigate the em-

pirical measure $\hat{\mu}_n := \frac{1}{n} \sum_i^n \delta_{x_i}$.

Lem. $\mathcal{Z} := \{ f \text{ is measurable. } |g| \le 1 \}$. For $\mu$

$\in \mathcal{J}$. s.t. $\mu = f \, dx$. $f \ge 0$. $f \in L^1$. Then:

$d_{\mathcal{Z}}(\mu \| \hat{\mu}_n) \equiv 2$. $\forall n \in \mathbb{N}$.

Pf: Set $g_n^w(x) = \sum_1^n \mathcal{I}_{\{X_j(w)\}}(x) - \mathcal{I}_{\{X_j(w), j \leq n\}^c}$

We have $\sup\limits_{\{g_n^w\} \cup \mathscr{wen}} d_g(\mu, \tilde{\mu}_n) = 2$

Df: $\mathcal{F}_{ac} := \{ g = \mathcal{I}_{(-\infty, \vec{a}]}, \vec{a} \in \mathbb{R}^k \}$. Glivenko-

Cantelli divergence $d_{ac}(\mu, \nu) = d_{g_{ac}}(\mu, \nu)$

$= \sup\limits_{\vec{a}} | F_\mu(a) - F_n(a) |$. $F_\mu$ is d.f. of $\mu$.

Rmk: i) By Lem. above, we restrict $\mathcal{F}$

on a small family.

ii) $d_{ac}$ is truly metric. Since:

$F_\mu(a) = F_\nu(a), \forall a \Rightarrow \mu \sim \nu$.

Thm. ( Glivenko - Cantelli)

$\mathcal{F} = M_b^1(\mathbb{R}^k)$ is learnable w.r.t. $d_{ac}$ by

SLA $\tilde{\mu}_n := \frac{1}{n} \sum_1^n \delta_{x_j}$

Rmk: It holds for $\mathbb{R}^k, \forall k \geq 1$.

Pf: Set $F_X(a-) = X_* \mathbb{P}(-\infty, a)$

$F_X(a) = X_* \mathbb{P}(-\infty, a]$.

and similarly for $\tilde{F}_n(a), \tilde{F}_n(a-)$.

Set $q_{\frac{j}{N}} = \inf \{x \mid F_X(x) \geq \frac{j}{N}\}$. $\frac{j}{N}$-quantile

of $X \sim \mu$. $q_0 = 0$. $q_1 = \infty$.

Let $\Gamma_n^{(N)} := \max_{1 \leq j \leq N-1} |F_X(q_{\frac{j}{N}}) - \widehat{F}_n(q_{\frac{j}{N}})| \vee$
$$|F_X(q_{\frac{j}{N}}-) - \widehat{F}_n(q_{\frac{j}{N}}-)|.$$

By SLLN. $\Gamma_n^{(N)} \to 0$ $(n \to \infty)$. $\mathbb{P}$-a.s.

For $\forall x \in \mathbb{R}^1$. $\exists j$. s.t. $x \in (q_{\frac{j-1}{N}}, q_{\frac{j}{N}})$.

So: $\widehat{F}_n(x) \leq \widehat{F}_n(q_{\frac{j}{N}}-) \leq F_X(q_{\frac{j}{N}}-) + \Gamma_n^{(N)}$
$$\leq F_X(x) + \frac{1}{N} + \Gamma_n^{(N)}.$$

$(N.tc : F_X(q_{\frac{j-1}{N}}-) - F_X(q_{\frac{j}{N}}-) \leq \frac{1}{N})$.

Similar for $\widehat{F}_n(x) \geq F_X(x) - \frac{1}{N} - \Gamma_n^{(N)}$.

So. $d_{dc}(\widehat{\mu}_n, \mu_X) \leq \frac{1}{N} + \Gamma_n^{(N)} \xrightarrow[n \to \infty]{N \to \infty} 0$. $\mathbb{P}$-a.s.

Rmk: For $J = M_+^1(\mathbb{R}^d)$, it's only learnable for

some weak distance like $d_{dc}$.