# Div. from Info. Theory

(1) Motivation from odds:

If we want to decide whether given data is from dist. $\mu$ or $\nu$:

i) If $\exists A, B \in \mathcal{B}_{\mathbb{R}^d}$. st. $\mu(A) > 0$, $\nu(B) > 0$ and $\mu(B) = \nu(A) = 0$. Assume sample $X = (x_1 \cdots x_n)$ is from $\mu$ of size $n$. So with proba $= 1 - (1 - \mu(A))^n$ we can observe at least one sample is from $A$. $\Rightarrow$ Make decision that $X$ is from dist. $\mu$.

Rmk: $1 - (1 - \mu(A))^n \xrightarrow{n \to \infty} 1$. So the conclusion can be reached exponentially fast.

ii) If $\mu \sim \nu$. (i.e. $\mu \ll \nu$ & $\nu \ll \mu$). Then: the decision become uncertained no matter how many samples we know.

Thm. (Radon-Nikodym)

For $\mu, \nu \in \mathcal{M}_1^+(\mathbb{R}^d)$. $\nu \ll \mu$. Then $\exists Z(x) \in$

$L^1(d\mu)$, s.t. $Z(x) \geq 0$ $\mu$-a.s. $dV(x) = Z(x) d\mu(x)$

**Pf:** Consider $L^2(\frac{1}{2}(\mu+V))$. We see:

$$|L_\mu(g)| \leq 2\|g\|_{L^2(\frac{1}{2}(\mu+V))}. \text{ So } L_\mu(\cdot) \text{ is}$$

BLO on $L^2(\frac{1}{2}(\mu+V)) \subseteq L^2(\mu) \cap L^2(V)$.

By Riesz Thm. $\exists \, s \in L^2(\frac{1}{2}(\mu+V))$, s.t.

$$\int g \, d\mu = \int g s \, d\tfrac{1}{2}(\mu+V). \quad \forall g \in L^2(\tfrac{1}{2}(\mu+V))$$

$$\Rightarrow \int g(2-s) \, d\mu = \int g s \, dV.$$

Set $A = \{s \leq 0\}$. We see that:

$$0 \leq \mu(A) \leq \tfrac{1}{2} \int I_A (2-s) \, d\mu = \int I_A s \, dV$$

$$\leq 0 \quad \Rightarrow \quad \mu(A) = 0 = V(A).$$

So: $V(B) = \int I_B \cdot \dfrac{s(x)}{s(x)} \, dV$

$$= \int I_B \, \frac{2-s}{s} \, d\mu. \quad \Rightarrow \quad Z = \frac{2-s}{s}$$

Let $A = \{Z < 0\}$. Then:

$$0 \leq V(A) = \int I_A \underset{(+)}{Z} \underset{(+)}{d\mu} \leq 0. \quad \Rightarrow \quad \mu(A) = 0.$$

**Rmk:** $Z(x) = dV/d\mu (x)$ is bdd of $V$ v.s.

$\mu$ at observed data pt. $X$.

(2) **Kullback - Leibler Div.:**

**Def**: For $\mu, \nu \in \mu_1^+(\mathbb{R}^d)$, $\nu \ll \mu$. KL div. is

$$d_{KL}(\mu \| \nu) = \int (\log \frac{d\mu}{d\nu}(x)) \cdot d\mu(x)$$

**Rmk**: i) If $\nu \not\ll \mu$. We set $d_{KL}(\mu \| \nu) = \infty$

ii) If $\mu \sim \nu$. We have $d_{KL}(\mu \| \nu)$

$$= \int -\log(\frac{d\nu(x)}{d\mu(x)}) \, d\mu(x).$$

**Lem**. (Jensen inequi.)

$q : \underset{interval}{I \subset \mathbb{R}^1} \to \mathbb{R}^1$. Convex. For $X$. $q(x) \in L^1$.

$$\Rightarrow \mathbb{E}(q(X)) \geq q(\mathbb{E}(X))$$

Besides if $q$ is strictly convex. and $X$

isn't deterministic. $\Rightarrow \mathbb{E}(q(X)) > q(\mathbb{E}(X))$

**Thm**. For $\mu, \nu \in \mu_1^+(\mathbb{R}^d)$. Then: $d_{KL}(\mu \| \nu) \geq 0$ and

$$d_{KL}(\mu \| \nu) = 0 \iff \mu = \nu.$$

**Rmk**: $d_{KL}$ is a divergence but not a

metric on $\mu_1^+(\mathbb{R}^d)$.

Claim: There's no metric $d(\cdot, \cdot)$ on

$\mu_1^+(\mathbb{R}^d)$. s.t. $d \sim d_{KL}$.

Pf: It's because $d_{KL}$ isn't sym:

$$d_{KL}(\mathcal{U}[0,\lambda m] \| \mathcal{U}[0,1]) \xrightarrow{m \to \infty} 0. \quad (\lambda m \downarrow 1)$$

But $d_{KL}(\mathcal{U}[0,1] \| \mathcal{U}[0,\lambda m]) \equiv \infty$

Since $\frac{d\mathcal{U}[0,\lambda m]}{d\mathcal{U}[0,1]} = \frac{1}{\lambda m}$ and

$\mathcal{U}[0,\lambda m] \not\ll \mathcal{U}[0,1]$.

$\underline{Pf}$: $d_{KL}(\mu \| \nu) = \int \left( \log \frac{d\mu}{d\nu} \right) \frac{d\mu}{d\nu} d\nu$

$f(t) = t \log t$ is strictly convex on $\mathbb{R}^{>0}$.

$\Rightarrow \mathbb{E}_\nu(f(x)) \geq f(\overline{\mathbb{E}_\nu(x)}) = 0.$

with "$=$" holds if $X = \frac{d\mu}{d\nu} \equiv const.$

( But $\mu(\mathcal{X}^d) = const. \cdot \nu(\mathcal{X}^d) = 1 \Rightarrow c = 1$ )

$\underline{Def}$: For $\mu, \nu \in \mathcal{M}_1^+(\mathcal{X}^d), \mu \sim \nu.$ Jensen-Shannon

div. is $d_{JS}(\mu \| \nu) = \frac{1}{2}(d_{KL}(\mu \| \frac{\mu + \nu}{2}) + d_{KL}(\nu \| \frac{\mu + \nu}{2}))$

$\underline{Rmk}$: $d_{JS}$ is also a div. which's of

importance in GANs.

# Comparison of Topo.

Next, we want to prove:

$$d_{KL} \Rightarrow d_{TV} \Rightarrow d_{KR} \overset{cpt}{\underset{set}{\Rightarrow}} d_{W_{1}CS_{1}} \Rightarrow \text{weak topo}$$

$$d_{LS} \hookrightarrow$$

$$\underbrace{\qquad\qquad}$$

equiv. in cpt set.

Rmk: These relation can help us transit some property (e.g. learnability) from one topo. to another topo.

Prop. (Gibbs' variational principle)

$X: (\mathbb{R}^d, \mathcal{B}_{\mathbb{R}^d}, \mu) \longrightarrow (\mathbb{R}, \mathcal{B}_{\mathbb{R}})$ with $(hF \; \gamma_{x,\mu}(q)$

Then: $\forall q > 0$. Set $\overline{E}_\nu(X) = -\infty$. if $\overline{E}_\nu(X_-) = -\infty$.

$$\gamma_{x,\mu}(q) = \sup_{\nu \in \mathcal{M}_1(\mathbb{R}^d)} (q \overline{E}_\nu(X) - d_{KL}(\nu \| \mu))$$

Rmk: Note that we don't restrict on $L'-$r.v. on condition.

Pf: Set $X_n = X \wedge n$. And define:

$$\frac{d\mu_{q,n}(X)}{} = e^{q X_n} d\mu(x) / \overline{E}_\mu(e^{q X_n})$$

( It's well-def since $X_n$ is upper-bdd)

For $k_{KL}(v\|p_t) < \infty$. (Since $k_{KL}(v\|\nu) = 0$

and RHS take supreme $\geq 0$)

$\gamma_{X_n, \mu}(t) - k_{KL}(v\|\mu_{t,n}) =$

$\gamma_{X_n, \mu}(t) + \int \log(d\mu_{t,n}/dv)\, dv =$

$\gamma_{X_n, \mu}(t) + \int \{\log(d\mu_{t,n}/d\mu) + \log(d\mu/dv)\}\, dv$

$= \gamma_{X_n, \mu}(t) + E_v(tX_n) - \gamma_{X_n, \mu}(t) - k_{KL}(v\|\mu).$

$= E_v(tX_n) - k_{KL}(v\|\mu).$

Take $\sup_{v \in \mu_{t,n}^+(t)}$ on both sides:

LHS $= \gamma_{X_n, \mu}(t)$ by: $v = \mu_{t,n}$ is optimal.

Next, we take $\sup_n$ on both sides:

By MCT. $\gamma_{X_n, \mu}(t) \uparrow \gamma_{X, \mu}(t)$ with

$$E_v(tX_n) \uparrow E_v(tX).$$

( It's consistent with $\overline{E}_v(X_-) = -\infty$

since $(X_n)_- = X_- \Rightarrow \overline{E}_v(X_n) = -\infty$)

<u>Lem.</u> $\sup_x \sup_g f(x,g) \overset{a)}{=} \sup_g \sup_x f(x,g) \overset{b)}{=} \sup_{x,g} f(x,g).$

$\underline{Pf:}$ Note $\sup_g f(x,g) \geq f(x,g). \Rightarrow a) : \geq \checkmark$

And by symm. $\Rightarrow$ a) : $\leq$ holds as well.

For b). Note: $\exists (x_k, y_k).$ s.t.

$f(x_k, y_k) \to \sup\limits_{x,y} f(x,y)$. So we have:

$$\sup\limits_{x} \sup\limits_{y} f(x,y) \geq \sup\limits_{y} f(x_k, y) \geq f(x_k, y_k)$$

Also $f(x,y) \leq \sup\limits_{x,y} f(x,y) \Rightarrow "\leq"$ holds.

$$\Rightarrow \Psi_{X,\mu}(\tau) = \sup\limits_{\sim} \sup\limits_{v} \left( \overline{\mathbb{E}}_v(\tau X_n) - d_{KL}(v \| \mu) \right)$$

$$= \sup\limits_{v} \sup\limits_{\sim} \overline{\mathbb{E}}_v(\tau X_n) - d_{KL}(v \| \mu)$$

$$= \sup\limits_{v \in \mathcal{P}_2(\mathbb{R}^d)} \overline{\mathbb{E}}_v(\tau X) - d_{KL}(v \| \mu).$$

Thm. (Bobkov - Götze)

$\mu \in \mathcal{P}_2(\mathbb{R}^d)$. $\delta$ is metric on $\mathbb{R}^d$. s.t. $x \in \mathbb{R}^d$

$\mapsto \delta(x,y)$ is Borel - measurable. $\forall y \in \mathbb{R}^d$. For

i) For $X \sim \mu$, $g \in Lip(\delta)$. $Y_g = g(X) - \mathbb{E}_\mu(g(x))$

is subgaussian with var. proxy $\sigma^2 > 0$

that doesn't depend on $g$.

ii) $k_{W_1(\delta)}(v \| \mu) \leq (2\sigma^2 d_{KL}(v \| \mu))^{\frac{1}{2}}$. $\forall v \in \mathcal{P}_2^+$.

Then: We have i) $\iff$ ii).

Pf: i) $\implies \Psi_{Y_g, \mu}(\tau) \leq \frac{1}{2}\sigma^2 \tau^2$. $\forall g \in Lip(\delta)$. $\forall \tau \geq 0$

$$\Longrightarrow \sup_{g \in Lip(\leq\delta)} \sup_{\alpha > 0} \sup_{v \in M^1} \left( \bar{E}_{v}(\alpha Y_g) - \lambda_{KL}(v \| \mu) - \tfrac{1}{2}\sigma^2\alpha^2 \right)$$

$$\leq 0 \quad \text{from variational principle.}$$

$$LHS \overset{Y_g = \cdots}{=} \sup_{v} \sup_{\alpha} \sup_{g} \left\{ \alpha \left( \int g \, dv - \int g \, d\mu \right) - \frac{\sigma^2}{2}\alpha^2 - \lambda_{KL}(v \| \mu) \right\}$$

$$\overset{\sup_{g \in Lip \leq \delta}}{=} \sup_{v} \sup_{\alpha} \left( \alpha \, d_{W, (\leq\delta)}(v \| \mu) - \tfrac{1}{2}\sigma^2\alpha^2 - \lambda_{KL}(v \| \mu) \right)$$

$$\overset{\alpha = \square}{=} \sup_{v} \left\{ (2\sigma^2)^{-1} d_{W, (\leq\delta)}(v \| \mu)^2 - \lambda_{KL}(v \| \mu) \right\}$$

So: $LHS \leq 0 \iff$ ii) holds.

e.g. Consider $\delta = 1.1$. $B \subseteq \mathbb{R}^d$. bdd domain with radius $r_B$, i.e. $\exists \, g_\varepsilon \in B$. s.t. $\sup_{x \in B} |g_\varepsilon - x|$

$\leq r_B + \varepsilon$. Next, restrict $v, \mu \in M^1_i(B)$:

First note Thm above works when we replace $Lip \leq \delta$ by $Lip_2 \leq \delta$, and $Y_g = Y_{g^*}$

if $g - g^* = const$. So let $g = g_\varepsilon$. Then:

$|g^*(x)| \leq r_B + \varepsilon$, $v$-a.s. $\forall v \in M^1_i(B)$. $g \in Lip_{g_2} \leq \delta$.

$\Rightarrow$ By Hoeffding inequi.: $\sigma^2 = (r_B + \varepsilon)^2$.

So: $d_{W, (\leq\delta)}(v \| \mu) \leq (r_B + \varepsilon)(2\lambda_{KL}(v \| \mu))^{\frac{1}{2}}$

$$\overset{\varepsilon \to 0}{\longrightarrow} r_B (2 \lambda_{KL}(v \| \mu))^{\frac{1}{2}}.$$

Note that Bobkov & Götze Thm above can be applied in any metric $\delta$ on $\mathscr{X}^d$.

**Def**: metric of the discrete topo $\delta_L(x,y)$ is:

$$\delta_L(x,y) = \mathbb{1}_{\{x \neq y\}}. \quad \forall x,y \in \mathscr{X}^d.$$

(Rmk: It's real metric. (check triangle inequ.))

**Lem.** i) $M_1^{+,\delta_L}(\mathscr{X}^d) = M_1^+(\mathscr{X}^d)$.

ii) $Lip(\delta_L) = \{g \in L^\infty(\mathscr{X}^d) \mid \sup_{x^d} g - \inf_{x^d} g \leq 1\}$.

iii) $d_{TV}(\mu,\nu) = 2\, d_{W,(\delta_L)}(\mu,\nu)$ on $M_1^+(\mathscr{X}^d)$.

**Pf**: i) Since $\{x\} \in B_{\mathscr{X}^d} \Rightarrow y \longmapsto \delta_L(y,x)$ is

measurable for $\forall x$.

And note: $\int \delta_L(x,y)\, d\nu = 1 - \nu(\{x\}) \leq 1$.

ii) "$\leq$": $\sup g - \inf g = \sup_{x,y} |g(x) - g(y)|$

$$\leq \delta_L(x,y) \leq 1.$$

"$\supseteq$": $\sup_{x,y} |g(x) - g(y)| \leq$

$$(\sup g - \inf g)\, \delta_L(x,y) \leq \delta_L(x,y).$$

iii) Set $g^*(x) = g(x) - \frac{1}{2}(\sup g + \inf g)$ for

$g \in Lip(\delta_L) \overset{ii)}{\Rightarrow} \|g^*\|_\infty \leq \frac{1}{2}. \quad g^* \in Lip(\delta_L)$

Note $\quad k_{TV}(\mu,\nu) = 2 \sup\limits_{g\in L^\infty, \|g\|_\infty \le \frac{1}{2}} |\int g\,d\mu - \int g\,d\nu|$

$$\le 2 \sup\limits_{g\in Lip(d_\lambda)} |\int g\,d\mu - \int g\,d\nu|.$$

$$= 2\, k_{W_1(d_\lambda)}(\mu,\nu)$$

$$= 2 \sup\limits_{g\in Lip(d_\lambda)} |\int g^*\,d\mu - \int g^*\,d\nu|.$$

$$\le 2 \sup\limits_{g^*\in L^\infty, \|g^*\|_\infty \le \frac{1}{2}} |\int g^*\,d\mu - \int g^*\,d\nu| = k_{TV}(\mu,\nu).$$

Thm. (Pinsker's Inequality)

$$k_{TV}(\mu,\nu) \le (2\, k_{KL}(\nu\|\mu))^{\frac{1}{2}} \quad \text{on } M_1^+(\mathfrak{X}).$$

Cor. $KL$-topo is stronger than $TV$-topo

and Radon-topo.

Pf: By Lem ii) above. $\forall g\in Lip(d_\lambda)$. $\overset{Doeblin}{\Longrightarrow}$

$Y_g$ is subgaussian with var. proxy $1/4$.

$$S_1: \quad k_{TV}(\mu,\nu) \overset{Lem\ ii)}{=} 2\, k_{W_1(d_\lambda)}(\mu,\nu)$$

$$\overset{Thm}{\le} 2\left(2\cdot\frac{1}{4}\cdot k_{KL}(\mu,\nu)\right)^{\frac{1}{2}}.$$