

# $\mathcal{I}$ -Divergence

We want to learn about the entire hist.  $\mu$  rather expected value. Because exp. value comes with risk and in order to quantify the value w.r.t the hist. Approx. knowledge of hist.  $\mu$  is needed.

Next, we will consider  $\mathcal{I}$ -divergence.

Remark: For  $\mathcal{I}$  larger,  $\delta_{\mathcal{I}}$  obtains more guarantees for approx. accuracy.

(1) Topo & Metric:

Def: i) V.W. n.v.s.  $L = (V, \|\cdot\|_V) \rightarrow (W, \|\cdot\|_W)$   
is strongly conti. if it's BLD.

ii)  $(V', \|\cdot\|'_V)$  denote dual space of  
 $(V, \|\cdot\|_V)$ .

iii)  $\|g\|_\infty := \sup_{x \in \mathbb{R}^k} |g(x)|$ . Set  $\tilde{L}_B^* \subset \mathbb{R}^k$   
:= { $f$  measurable |  $\|f\|_\infty < \infty$ }.

Rank:  $M_+^t(\mathbb{R}^k)$  can be embedded into

$\langle L_\infty \rangle^\perp$  by set  $L_V(g) = \int g dV$

for  $V \in M_+^t(\mathbb{R}^k)$ .  $|L_V(g)| \leq \|g\|_\infty$ .

iv)  $\mu(\mathbb{R}^k)$  is finite signed measure on  $(\mathbb{R}^k, \mathcal{B}_{\mathbb{R}^k})$ . It can be embedded in  $\langle L_B^- \rangle^\perp$  as well: for  $\gamma = \alpha\mu + \beta V$ .

$$\begin{aligned} L_\gamma(g) &= \alpha L_\mu(g) + \beta L_V(g) \\ &= \alpha \int f d\mu + \beta \int f dV = \int g d\gamma. \end{aligned}$$

v) Equip  $\mu(\mathbb{R}^k)$  with norm  $\|\gamma\|_{TV} := \sup \{L_\gamma(f) \mid f \in L_B^+(\mathbb{R}^k), \|f\|_\infty = 1\}$ .

Rank: By Jordan decom.,  $\gamma = \gamma^+ - \gamma^-$ .

sz.  $\gamma^+ \perp \gamma^-$ . So  $\exists A^\pm \in \mathcal{B}_{\mathbb{R}^k}$ . Sz.

$f = \bar{f}_{A^+} - \bar{f}_{A^-}$  attains its sup.

vi)  $\mathcal{F}_{TV} := L_B^{(\infty)}(\mathbb{R}^k)$ .  $\|\gamma_{TV}(\mu, \nu)\| = \|q_{TV}(\mu, \nu)\|$   
 $= \|\mu - \nu\|_{TV}$  for  $\mu, \nu \in \mu_+^t(\mathbb{R}^k)$ .

Rank: i) Set  $\mathcal{F}_{TV} := \mathcal{C}_B(\mathbb{R}^k)$ . And we restrict on it. We can obtain

Rakon-Norm:  $\|\cdot\|_{R0}$ . Which's weaker than  $\|\cdot\|_{TV}$ . Also we can also consider  $g_c := C_c \langle \cdot \rangle$ .

ii) We can also introduce  $d_{TV}(\mu, \nu) = 2 \sup_{A \in \mathcal{B}(X)} |\mu(A) - \nu(A)|$ . i.e. set  $g = \epsilon^{-2} I_A (A \in \mathcal{B}(X))$ .

Lem.  $X_n \sim \mu_n$ ,  $X \sim \mu$ .  $\epsilon p^k - r.v.$ . Then: We have

$X_n \xrightarrow{\text{P}} X \Leftrightarrow \mu_n \rightarrow \mu$  in  $C_c(C^{R^k})'$ .

Pf: prove that  $\mu_n \rightarrow \mu$  in  $(C_c(C^{R^k}))'$  ( $\Leftrightarrow$ )  
 $\mu_n \rightarrow \mu$  in  $(C_b(C^{R^k}))'$ .

For  $g \in (C_b(C^{R^k}))'$ .  $B = \|g\|_\infty$ . Choose  $k$

so.  $\mu \in B_k(0)^\perp \leq \epsilon/B$ .

$\varphi$  is bump func. st.  $\varphi = 1$  on  $B_k(0)$ .

$\varphi \geq 0$ .  $\|\varphi\|_\infty = 1$ .  $\varphi \in C_c^\infty(C^{R^k})$ .

$$\int g d\mu_n = \int \varphi g d\mu_n + \int (1-\varphi) g d\mu_n$$

$$|\int (1-\varphi) g d\mu_n - \int (1-\varphi) g d\mu| \leq$$

$$B(1 - \int \varphi d\mu_n) + B \leq 1 - \int \varphi d\mu \rightarrow$$

$$2B(1 - \int \varphi d\mu) \leq B \cdot \frac{\epsilon}{B} = \epsilon \rightarrow 0.$$

Cor. For  $\mathcal{T} = \mu^+ \cap \mathbb{R}^k$ , to be the space of p.m.s < $\ll$   $f_x$  or ' $\mu^k$ '. i.e.  $\mu(f_x) = f(x)x$ . Then:  $\forall \mu, \nu \in \mathcal{T}$ . we have:

$$k_{\tau\nu}(\mu, \nu) = k_{C_C}(\mu, \nu) = k_{\text{LD}}(\mu, \nu).$$

Pf.: For the latter part. we set

$\varphi_n$  is bump func. on  $B_{\mathbb{R}^k}^{\text{cusp}}$ .

and  $\|\varphi_n\|_\infty \leq 1$

And use DCT. we can get:

$$\int \varphi_n g dm \rightarrow \int g dm.$$

Rmk: Note that  $C_B, L_B^\infty$  are not

separable. So it may lead to some measurability problem when considering SLAs.

(2) Wasserstein metric:

Df: For  $\delta$  metric on  $\mathbb{R}^k$ . s.t.  $\delta(x-y) \in B_{\mu, \nu, \mathbb{R}^k}$ .  $\forall x, y$ .

$\Rightarrow \text{Lip}(\delta) = \{f: \mathbb{R}^k \rightarrow \mathbb{R}, f \text{ has Lip-const} = 1\}$ .  $\text{Lip}_1(\delta) := \text{Lip}(\delta) \cap \{f' | f'(y) = 0\}$

Remark:  $g = x \in \text{Lip}_0^{\delta}$ .  $\mu \sim \text{Cauchy}(1)$ ,  $L_{\mu}(g) = \infty$   
 So  $L_{\mu}(g)$  doesn't have to exist.

ii)  $\mu^{\delta \times \mathbb{R}^k} := \{ \nu \in \mathcal{M}^{\delta \times \mathbb{R}^k} \mid \exists \gamma \in \mathbb{R}^k, \text{ s.t.}$

$$\int \delta(x, y) \lambda(\nu)(x) < \infty \}.$$

iii) Wasserstein 1-norm  $\|\mu - \nu\|_{W,cs} = \sup \{ L_{\mu-\nu}(g) \mid g \in \text{Lip}_{\mathbb{R}^k}^{\delta}, \text{ for some } \gamma \in \mathbb{R}^k$   
 and  $\mu, \nu \in \mathcal{M}^{\delta \times \mathbb{R}^k}$ .

iv) Set  $\mathcal{M}_+^{\delta, t \times \mathbb{R}^k} = \mathcal{M}_+^t \times \mathbb{R}^k \cap \mathcal{M}^{\delta \times \mathbb{R}^k}$ .

Lip metric  $\| \cdot \|_{W,cs} : = \|\mu - \nu\|_{W,cs}$   
 is defined on  $\mathcal{M}_+^{\delta, t \times \mathbb{R}^k}$ .

Lemma: For  $\mu, \nu \in \mathcal{M}^{\delta \times \mathbb{R}^k}$ . We have :

i)  $\|\mu\|_{W,cs} < \infty$ .

ii) If  $\mu \times \mathbb{R}^k = \nu \times \mathbb{R}^k$ . Then we can

replace  $\text{Lip}_{\mathbb{R}^k}^{\delta}$  by  $\text{Lip}^{\delta}$ . i.e.

$$\|\mu - \nu\|_{W,cs} = \sup_{\text{Lip}^{\delta}} \{ L_{\mu-\nu}(g) \mid$$

iii)  $\|\cdot\|_{W,cs}$  is a semi-norm on  $\mathcal{M}^{\delta \times \mathbb{R}^k}$ .

iv) If  $\delta$ -topo is weaker than  $t$ -topo.  
 initial  $\delta$ -topo.  $\subset X_n \xrightarrow{\delta} x \Rightarrow X_n \xrightarrow{t} x$ . Then :

$L_{W,cs}(\mu, \nu) = 0 \iff \mu = \nu$  for  $\mu, \nu \in \mathcal{M}_1^{s.t.}(\mathbb{R}^d)$ . So it's metric on  $\mathcal{M}_1^{s.t.}(\mathbb{R}^d)$ .

Pf: i) Note  $g(x) = |g(y)| + \delta(x,y) = \delta(x,y)$

ii) Let  $g^*(x) = g(x) - g(y)$ . Then:

$L_{\mu-\nu}(g) = L_{\mu-\nu}(g^*)$ , from add.

iii)  $\| \cdot \|_{W,cs}$  is supremum of seminorm.

iv)  $L_{W,cs}(\mu, \nu) = \sup_{\text{Lip-const.}} L_\mu(g) = L_\nu(g)$

f.r.  $\forall g$  is Lip-conti.

Next, we prove  $\mu(A) = \nu(A)$  for

$\forall A$  closed in  $\mathbb{R}^d$ .

Set  $J_{\varepsilon,A}(x) = (1 - \frac{\varepsilon}{\delta} \delta(x, A)) \vee 0$ .

$\Rightarrow J_{\varepsilon,A}$  is Lip-conti. with const.

$= \frac{1}{\varepsilon}$ . And  $J_{\varepsilon,A} \downarrow J_A$  because

$\delta(x, A) = 0 \stackrel{\text{const.}}{\Rightarrow} \inf_y |x-y| = 0$ .

Applying PCT on  $\int J_{\varepsilon,A} d\mu = \int J_{\varepsilon,A} d\nu$ .

Def.  $\delta(x, y) = |x-y|$  and  $B \subset \mathbb{R}^d$  bdd. Let

$M_1^{s.t.}(B) = \sum \mu \in \mathcal{M}_1^{s.t.}(\mathbb{R}^d) \mid \text{supp}(\mu) \subset B \}$ .

Let  $r_B = \inf_{g \in \mathcal{B}} \sup_{x \in B} \delta(x, g)$ . Then: we have.

$L_{W_{(S)}}(\mu, \nu) \leq r_B L_{RD}(\mu, \nu)$  for  $\forall \mu, \nu \in \mathcal{M}_1^+(B)$ .

Pf:  $\exists g_2 \in \mathcal{B}$ . s.t.  $\sup_{x \in B} \delta(x, g_2) \leq \varepsilon + r_B$ .

S. :  $\forall g \in \text{Lip}_{g_2}(\delta)$ .  $\Rightarrow g = r_B + \varepsilon$  on  $B$ .

Let  $\bar{g} = g(x) I_B + (r_B + \varepsilon) I_{B^c}$ .

$$\Rightarrow L_{\mu-\nu}(g) = L_{\mu-\nu}(\bar{g}).$$

$$= (r_B + \varepsilon) L_{\mu-\nu}\left(\frac{\bar{g}}{r_B + \varepsilon}\right)$$

$$\leq (r_B + \varepsilon) \| \mu - \nu \|_{RD}.$$

C'. Under condition above. we have:

i) If  $B$  is open. then  $\| \cdot \|_{W_{(S)}}$  is weaker than  $\| \cdot \|_{RD}$  or  $\| \cdot \|_{TV}$ .

ii)  $\| \cdot \|_{W_{(S)}} - \text{topo} = \text{weak convergence}$

(3) Metrics w.r.t density / para. :

① For  $\lambda \mu = f_\mu \lambda x$ ,  $\lambda \nu = f_\nu \lambda x$ . We set:

$$\lambda \rho(\mu, \nu) = \| f_\mu - f_\nu \|_{L^1} \text{ on } \mathcal{I} = \sum \mu \in \mathcal{M}_1^+(\mu < \lambda x)$$

Rank: For  $p=1$ . Note  $\text{sgn}(f_1 - f_2) \in L^{\frac{1}{2}}$ .

So it corresponds to  $\| \cdot \|_{TV} = \| \cdot \|_{K_0}$ .

For  $p > 1$ . They're not comparable.

② For  $\mathbb{Q} \subset \mathbb{R}^k$ . Consider map  $\mu_{(.)} : \mathbb{Q} \rightarrow \mathbb{M}^+(\mathbb{R}^k)$ ,

is injective. And for  $v, v' \in \text{Im}(\mu)$ . Set

$$d_{\mathbb{Q}}(v, v') = |\theta - \theta'| \text{ where } v = \mu_\theta, v' = \mu_{\theta'}.$$

Rank: It can develop para. stat. approach

$$\kappa(\hat{\mu}_n \| \mu) = \kappa(\mu_{\theta_n} \| \mu_\theta) \leq C |\theta_n - \theta| \rightarrow 0.$$

Lem.  $d_{\mathbb{Q}}(., .)$  is divergence on  $\mathbb{M}^+(\mathbb{R}^k)$  and

assume  $\mu_{(.)} : \mathbb{Q} \rightarrow \mathbb{M}^+(\mathbb{R}^k)$  is injective

and conti. w.r.t.  $\lambda_0 = 1 \cdot 1$ . Then:

i)  $\kappa_{\mathbb{Q}}$ -topo is stronger than  $\kappa$ -topo.

ii) If  $\mathbb{Q}$  is cpt. Then:  $\lambda_0 \sim \lambda$  on  $\text{Im}(\mu_{(.)})$

Pf: i)  $\theta_n \xrightarrow{1 \cdot 1} \theta_0 \stackrel{\text{conti.}}{\Rightarrow} \kappa(\mu_{\theta_0} \| \mu_{\theta_n}) \rightarrow 0$ .

ii) If  $\kappa(\mu_{\theta_0} \| \mu_{\theta_n}) \rightarrow 0$ .  $|\theta_n - \theta_0| \rightarrow 0$

$$\Rightarrow \exists \epsilon > 0, \forall n_k. (\epsilon \cdot (\theta_{n_k}) \subset \mathbb{Q}_\epsilon = \mathbb{Q} / \beta_\epsilon(\theta_0))$$

Note  $\mathbb{W}_\varepsilon$  is also cpt.

So  $\delta = \inf_{\theta \in \mathbb{W}_\varepsilon} d(\mu_{\theta_0}, \mu_\theta) > 0$  follows

from the univ. condition of  $\mu_{\cdot \cdot}$

(Otherwise it will contradict with  
that " $\mu_{\cdot \cdot}$  is injective and  $\lambda$   
is separating. i.e.  $\exists \theta_1 \neq \theta_2. \mu_{\theta_0} = \mu_{\theta^*}$ )

$\Rightarrow d(\mu_{\theta_0}, \mu_{\theta_k}) \geq \delta > 0$ .

which contradict with  $\mu_{\theta_k} \xrightarrow{\lambda} \mu_{\theta_0}$ .