

## L5: Two-sample problem and MMD

Two-sample problem

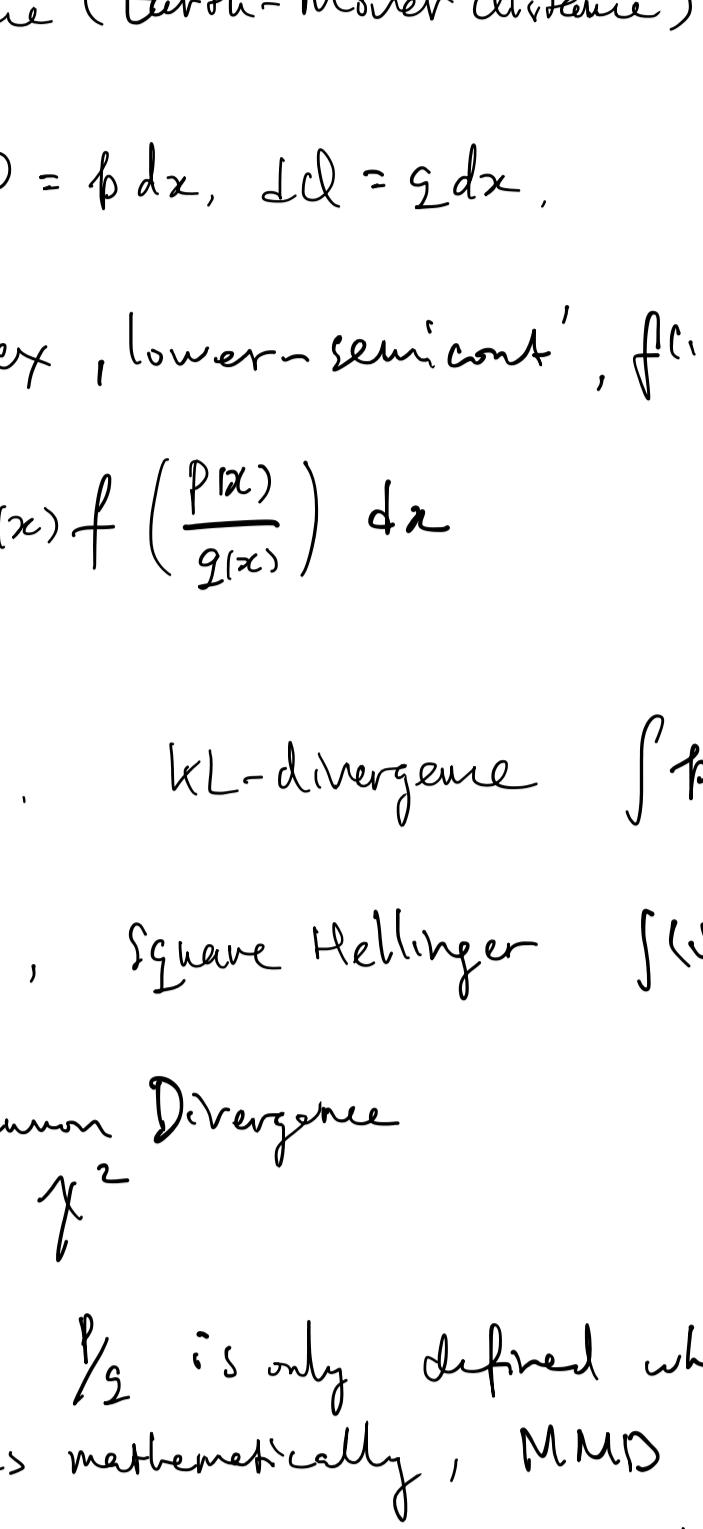
$$\text{given } X = \{x_i\}_{i=1}^{n_x}, Y = \{y_i\}_{i=1}^{n_y}, \\ x_i \sim p, y_i \sim q, i.i.d., X \perp\!\!\!\perp Y$$

Can we tell if  $p = q$  from the data?  
When  $p \neq q$ , where do they differ?

RK: A problem in the GAN loss. Related to the estimation of distance/divergence b/w two (unknown) distributions from finite sample.

Standard procedure:

- Compute test statistic  $T(X, Y)$
- Specify a threshold value  $\tau$
- Accept  $H_0$  if  $T(X, Y) < \tau$ , reject otherwise



In 1D: Kolmogorov-Smirnov, Wald-Wolfowitz runs, ...

We consider general  $P$  and  $Q$  in  $\mathbb{R}^d$ .

- MMD (maximum mean discrepancy)

Integral Probability Metrics (IPM)

$$D_{\text{IPM}}(P, Q) = \sup_{f \in \mathcal{H}_k} |\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{y \sim Q} f(y)|$$

$$\text{Eq } \mathcal{H}_k = \{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \leq 1\} \text{ kernel MMD},$$

$\mathbb{E}_f$   $\int p = \int p dx, \int q = \int q dx$ ,  $p, q$  densities

$$W_1(p, q) = \sup_{\substack{f: \mathbb{R} \rightarrow \mathbb{R} \\ \|f\|_{\text{Lip}} \leq 1}} |\int f(p-q)|$$

$$\text{Lip}(f) := \sup_{\substack{x \neq y \in \mathbb{R} \\ \|x-y\|_2 \leq 1}} \frac{|f(x) - f(y)|}{\|x-y\|_2} \text{ Lipschitz semi-norm}$$

Wasserstein-1 distance (Earth-Mover distance)

-  $f$ -divergence:  $\int p = \int p dx, \int q = \int q dx$ .

$f: \mathbb{R}_+ \rightarrow \mathbb{R}$ , convex, lower-semicontinuous,  $f(0) = 0$

$$D_f(P||Q) = \int_{\mathbb{R}} q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

$\text{Eq } f(u) = \min_u u \ln u$ . KL-divergence  $\int p \ln \frac{p}{q}$

$\text{Eq } f(u) = (\sqrt{u}-1)^2$ , Square Hellinger  $\int (p-\sqrt{q})^2$

Other: Jensen-Shannon Divergence  
Pearson  $\chi^2$

RK: density ratio  $\frac{p}{q}$  is only defined where  $q(x) > 0$ . This mathematically, MMD can handle when  $\text{supp}(p)$  and  $\text{supp}(q)$  differ.

- MMD is a pseudometric on  $\mathcal{P}(\mathbb{R})$   
the space of prob. measures.

$$(i) D_{\text{IPM}}(P, P) = 0$$

$$(ii) D_{\text{IPM}}(P, Q) = D_{\text{IPM}}(Q, P)$$

$$(iii) \text{Triangle inequality } D_{\text{IPM}}(P, R) \leq D_{\text{IPM}}(P, Q) + D_{\text{IPM}}(Q, R)$$

However, generally,  $D_{\text{IPM}}(P, Q) = 0 \not\Rightarrow P = Q$

- RKHS-MMD (kernel MMD)

If  $k$  the RKHS associated with PSD kernel  $k$  on  $\mathbb{X}$ ,  $(\mathbb{X}, \mu)$  the measure space  
 $P$  and  $Q$  are prob. distributions on  $\mathbb{X}$ , have densities.

kernel MMD: (the absolute value in def can removed)

$$\text{MMD}(p, q) = \sup_{\substack{f \in \mathcal{H}_k, \\ \|f\|_{\mathcal{H}_k} \leq 1}} \int f(p-q)$$

def Mean embedding of  $P$  is  $\mu_p = \mathbb{E}_{y \sim p} k(x, y)$

lemma If  $\mathbb{E}_{x \sim p} \sqrt{k(x, z)} < \infty$ , then  $\mu_p \in \mathcal{H}_k$ .

If  $\mu_p(x)$  is well defined b/w  $\|\mu_p\|_{\mathcal{H}_k} \leq \sqrt{\mathbb{E}_{x \sim p} k(x, z) \mathbb{E}_{y \sim p} k(y, z)}$ ,  
thus  $k(x, y)$  is integrable under  $\mathbb{E}_{x, y}$ .

$\mu_p = \mathbb{E}_{y \sim p} \phi(y)$ , then

$$\|\mu_p\|_{\mathcal{H}_k}^2 = \langle \mu_p, \mu_p \rangle_{\mathcal{H}_k} = \langle \mathbb{E}_{y \sim p} \phi(y), \mathbb{E}_{y \sim p} \phi(y) \rangle_{\mathcal{H}_k}$$

$$\stackrel{(c-s)}{=} \mathbb{E}_{x \sim p, y \sim p} \frac{\langle \phi(x), \phi(y) \rangle_{\mathcal{H}_k}}{\sqrt{\mathbb{E}_{x \sim p} k(x, x) \mathbb{E}_{y \sim p} k(y, y)}} \leq \infty$$

Prop. For  $f \in \mathcal{H}_k$ ,  $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}_k}$ , then

$$(\mathbb{E}_{x \sim p} - \mathbb{E}_{x \sim q}) f(x) = \langle f, \mathbb{E}_{x \sim p} \phi(x) - \mathbb{E}_{x \sim q} \phi(x) \rangle_{\mathcal{H}_k}$$

$$= \langle f, \mu_p - \mu_q \rangle_{\mathcal{H}_k}$$

thus  $\sup_{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mu_p - \mu_q \rangle_{\mathcal{H}_k}$  is achieved

when  $f = \frac{\mu_p - \mu_q}{\|\mu_p - \mu_q\|_{\mathcal{H}_k}}$ , if  $\|\mu_p - \mu_q\|_{\mathcal{H}_k} > 0$ .

and then  $\text{MMD}(p, q) = \|\mu_p - \mu_q\|_{\mathcal{H}_k}$ .  
(when  $\|\mu_p - \mu_q\|_{\mathcal{H}_k} = 0$ ,  $\text{MMD} = 0$ )

Thus  $\text{MMD}(p, q)^2 = \langle \mu_p - \mu_q, \mu_p - \mu_q \rangle_{\mathcal{H}_k}$

$$= \langle (\mathbb{E}_{x \sim p} - \mathbb{E}_{x \sim q}) \phi(x), (\mathbb{E}_{y \sim p} - \mathbb{E}_{y \sim q}) \phi(y) \rangle_{\mathcal{H}_k}$$

$$= (\mathbb{E}_{x \sim p} - \mathbb{E}_{x \sim q}) (\mathbb{E}_{y \sim p} - \mathbb{E}_{y \sim q}) k(x, y)$$

RK: The pf shows that the maximiser equals  $\langle \mu_p - \mu_q, \mu_p - \mu_q \rangle_{\mathcal{H}_k}$  up to a normalising constant. By definition.

$$\langle \mu_p - \mu_q, \mu_p - \mu_q \rangle_{\mathcal{H}_k} = \int_{\mathbb{X}} k(x, y) (p-q)(y) dy = w(x)$$

witness function of kernel MMD.

$\square$  kernel  $k$ ,  $\text{MMD}_k(p, q) = 0 \Rightarrow p = q$

Let  $\mathbb{X}$  be a compact metric space,  $M_b(\mathbb{X})$  the set of finite signed measures on  $\mathbb{X}$ .

def Kernel  $k$  is universal iff for any  $\mu \in M_b(\mathbb{X})$ ,

$$\int_{\mathbb{X}} k(x, z) d\mu(z) = 0 \Rightarrow \mu = 0$$

Prop. The def is equivalent to that

$$\int_{\mathbb{X} \times \mathbb{X}} k(x, y) d\mu(x) d\mu(y) > 0, \forall \mu \in M_b(\mathbb{X}) \setminus \{0\}. \quad (*)$$

$$\text{where } f = \frac{\mu_p - \mu_q}{\|\mu_p - \mu_q\|_{\mathcal{H}_k}}, \text{ if } \|\mu_p - \mu_q\|_{\mathcal{H}_k} > 0.$$

and then  $\text{MMD}(p, q) = \|\mu_p - \mu_q\|_{\mathcal{H}_k}$ .  
(when  $\|\mu_p - \mu_q\|_{\mathcal{H}_k} = 0$ ,  $\text{MMD} = 0$ )

thus if  $\text{supp}(\hat{\mu}) = \mathbb{R}^d$ , then r.h.s = 0 only when

$\hat{\mu}(\mathbb{R}^d) = 0 \Rightarrow \mu = 0$ .  
This proves the  $\Rightarrow$  direction of the prop. the other direction  $\Leftarrow$  as  $\mathbb{R}^d$ .

$\square$  Universal kernels in  $\mathbb{R}^d$ :

$$e^{-\|x-y\|^2/\sigma^2} \text{ Gaussian kernel}$$

$$e^{-\|x-y\|_1} \text{ Laplace kernel}$$

- Estimation of MMD from data

$\text{MMD}(p, q)$  is called population MMD, given finite samples  $X = \{x_i\}_{i=1}^{n_x}, Y = \{y_i\}_{i=1}^{n_y}$ ,  $x_i \sim p$  i.i.d.,  $y_i \sim q$  i.i.d.,  $X \perp\!\!\!\perp Y$ .

How to estimate MMD from  $X$  and  $Y$ ?

Denote  $T = \text{MMD}(p, q)^2 = \int \int k(x, y) (p-q)(y) dx dy$

$$= \frac{1}{n_x} \sum_{i=1}^{n_x} \hat{x}_i - \frac{1}{n_y} \sum_{j=1}^{n_y} \hat{y}_j$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j)$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{x}_i \rangle - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$

$$= \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{x}_i) - \frac{1}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(\hat{x}_i, \hat{y}_j) \langle \hat{x}_i, \hat{y}_j \rangle$$
</div