# Lecture 20: MMD in Practice

last lecture we looked at the basic properties
of kernel mean embeddings of distributions

$$P(x) \Rightarrow \mu \longmapsto \mu_{\kappa} := \mathbb{E}_{\underline{x} \sim \mu} K(x, \cdot) \in H_{\kappa}$$
$$= \int K(\underline{\xi}, \cdot) \mu(d\underline{\xi})$$

& the associated Maximum mean discrepancy

$$MMD(\mu, \mu') := \| \mu_{\kappa} - \mu'_{\kappa} \|_{\kappa}$$

The main take aways were:

- If $\int \sqrt{K(\underline{x}, \underline{x})} \mu(d\underline{x}) < + \infty$ then $\mu_{\kappa}$ is the Riesz. rep. of the bdd. lin operator

$$\phi \in H_{\kappa}^{*}, \quad \phi(f) = \int f(\underline{x}) \mu(d\underline{x})$$
$$= \langle f, \mu_{\kappa} \rangle_{\kappa}$$

- MMD satisfies all but one property of a metric, ie, we may have in general

$$MMD(\mu, \mu') = 0 \quad \text{while} \quad \mu \neq \mu'$$

①

As an example consider the linear kernel
$$K(\underline{x}, \underline{z}) = \underline{x}^T \underline{z}$$ then we have

$$\mu_K(\underline{z}) = \mathbb{E}_{\underline{x} \sim \mu} \underline{x}^T \underline{z} = \underline{m}^T \underline{z} \qquad \underline{m}, \underline{m}' \text{ are mean}$$
$$\mu'_K(\underline{z}) = \mathbb{E}_{\underline{x} \sim \mu'} \underline{x}^T \underline{z} = (\underline{m}')^T \underline{z} \qquad \text{of } \mu, \mu'$$

& so, $MMD(\mu, \mu')^2 = \|\underline{m} - \underline{m}'\|^2$ so that
$MMD(\mu, \mu') = 0$ so long as $\mu, \mu'$ have the same mean!

&bull; We introduced the idea of a
**Characteristic kernel**, ie, a kernel $K$
such that $MMD(\mu, \mu') = 0 \Longleftrightarrow \mu = \mu'$
for all $\mu, \mu' \in \mathbb{P}_K(\mathcal{X})$.

eg: RBF kernel $K(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|^2}{2\gamma^2}\right)$

Fourier $K(\underline{x}, \underline{x}') = \exp(i\underline{x}^T\underline{x}')$

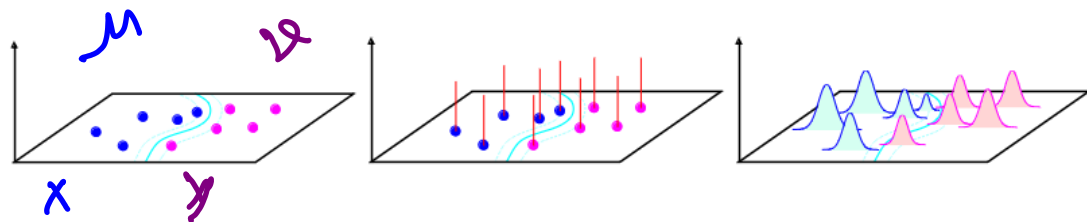In this lecture we want to consider more practical aspects of the MMD & kernel mean embeddings.

②

# 20.1 – empirical estimation of MMD

In most practical applications in statistics or data science we are given data sets

$$X = \{\underline{x}_1, \ldots, \underline{x}_n\} \subset \mathbb{R}^d$$
$$Y = \{\underline{y}_1, \ldots, \underline{y}_m\} \subset \mathbb{R}^d$$

& the general assumption is that these data sets are drawn (independently) from underlying measures $\mu, \nu$.



So, at best, assuming $\underline{x}_j$ & $\underline{y}_j$ are iid wrt $\mu, \nu$, we can approximate the MMD between the empirical distributions

$$\mu^n = \frac{1}{n} \sum_{j=1}^{n} \delta_{\underline{x}_j} \quad , \quad \nu^n = \frac{1}{m} \sum_{j=1}^{m} \delta_{\underline{y}_j}$$

It turns out, that the reproducing property allows us to write a very simple expression for computing the empirical MMD.

③

$$\mu^n_k = \int K(\underline{x}, \cdot) \mu^n (dx) = \frac{1}{n} \sum_{j=1}^{n} K(\underline{x}_j, \cdot)$$

$$\nu^m_k = \int K(\underline{x}, \cdot) \nu^m (dx) = \frac{1}{m} \sum_{j=1}^{m} K(\underline{y}_j, \cdot)$$

& so,

$$MMD(\mu^n, \nu^n)^2 = \| \mu^n_k - \nu^m_k \|^2_K$$

$$= \| \mu^n_K \|^2_K + \| \nu^m_k \|^2_k - 2 \langle \mu^n_k, \nu^m_k \rangle$$

$$= \left\| \frac{1}{n} \sum_{j=1}^{n} K(\underline{x}_j, \cdot) \right\|^2_k + \left\| \frac{1}{m} \sum_{j=1}^{m} K(\underline{y}_j, \cdot) \right\|^2_k$$

$$- \frac{2}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \langle K(\underline{x}_i, \cdot), K(\underline{y}_j, \cdot) \rangle$$

$$= \frac{1}{n^2} \sum_{i,j=1}^{n} K(\underline{x}_i, \underline{x}_j) + \frac{1}{m^2} \sum_{i,j=1}^{m} K(\underline{y}_i, \underline{y}_j)$$

$$- \frac{2}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} K(\underline{x}_i, \underline{y}_j)$$

$(\lambda)$

This convenient expression is one of the most attractive features of the MMD!

$(4)$

The question arises as to how good of an approx. is $MMD(\mu^n, x^m)$?

Thm: Let $k$ be a continues & PDS kernel on a separable metric space $X$ such that

$$\sup_{\underline{x} \in X} K(\underline{x}, \underline{x}) < C_k < +\infty.$$ Thm for any

$\delta \in (0,1)$ with Prob. at least $1-\delta$ we have

$$MMD(\mu, \mu^n) \leq \sqrt{\frac{C_k}{n}} + \sqrt{\frac{2C_k \, \ell_g \frac{1}{\delta}}{n}}$$

By triangle ineq. we have

$$MMD(\mu, x) \leq MMD(\mu, \mu^n) + MMD(\mu^n, x)$$

$$\leq MMD(\mu, \mu^n) + MMD(\mu^n, x^n) + MMD(x^n, x)$$

$$\Rightarrow |MMD(\mu, x) - MMD(\mu^n, x^n)| = G\left(\sqrt{\frac{1}{n}} \vee \sqrt{\frac{1}{m}}\right)$$

with high prob. ∴

⑤  In other words, convergence happens at monte carlo rate!

Computing MMD using the formula $\circledast$ is
in general an $O(n^2 d)$ operation (assuming $m = \Theta(n)$).
This can be prohibitive in some cases. Luckily
there are many work arounds, such as
random features.

Recall if we have a stationary kernel
$K(\underline{x}, \underline{x}') = K(\underline{x} - \underline{x}')$ then we could apprea
our kernel with

$$K^N(\underline{x}, \underline{x}') = \frac{1}{N} \sum_{j=1}^{N} \varphi_j(\underline{x}) \varphi_j(\underline{x}')$$

where $\varphi_j$ where the random features of our
kernel. We also showed that the RkHS of
$K^N$ consists of functions $f = \sum_{j=1}^{N} \frac{\alpha_j}{\sqrt{N}} \varphi_j$
with RKHS norm $\|f\|^2 = \|\underline{\alpha}\|_2^2$.

Based on this we can immediately see that

⑥

$$\text{MMD}_k(\mu, \nu)^2 \approx \text{MMD}_{k^N}(\mu, \nu)^2$$

$$= \| \mu_{k^N} - \nu_{k^N} \|_{k^N}^2$$

$$= \left\| \int k^N(\underline{x}, \cdot)\mu(d\underline{x}) - \int k^N(\underline{x}, \cdot)\nu(d\underline{x}) \right\|_{k^N}^2$$

$$= \left\| \int \sum_{j=1}^N \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \frac{1}{\sqrt{N}} \varphi_j(\cdot)\mu(d\underline{x}) \right.$$
$$\left. - \int \sum_{j=1}^N \frac{1}{\sqrt{N}} \varphi_j(\underline{x}) \frac{1}{\sqrt{N}} \varphi_j(\cdot)\mu(d\underline{x}) \right\|_{k^N}^2$$

$$= \left\| \sum_{j=1}^N \left( \int \frac{1}{\sqrt{N}} \varphi_j(\underline{x})\mu(d\underline{x}) - \int \frac{1}{\sqrt{N}} \varphi_j(\underline{x})\nu(d\underline{x}) \right) \frac{1}{\sqrt{N}} \varphi_j \right\|_{k^N}^2$$

$$= \frac{1}{N} \sum_{j=1}^N \left| \int \varphi_j(\underline{x})\mu(d\underline{x}) - \int \varphi_j(\underline{x})\nu(d\underline{x}) \right|^2$$

$$\approx \frac{1}{N} \sum_{j=1}^N \left| \frac{1}{n} \sum_{\ell=1}^n \varphi_j(\underline{x}_\ell) - \frac{1}{m} \sum_{\ell=1}^m \varphi_j(\underline{y}_\ell) \right|^2$$

for $\underline{x}_\ell \sim \mu,\ \underline{y}_\ell \sim \nu$.

⑦

Observe that this formula not only involves the mean of $N$ random features. It is also highly parallelizable since the $\varphi_j(x_i)$ can be computed independently for each feature.

Also observe that the above calculation can be repeated foor any other low-rank approx. to our kernel.

$$K(\underline{x}, \underline{x}') \approx \sum_{j=1}^{N} \psi_j(\underline{x}) \psi_j(\underline{x}')$$