# 8 - Kernel Methods: Properties and Applications

Feature map $\phi$

Kernel $k(\cdot, \cdot)$

unique

RKHS $\mathcal{F}$

Recall: 1 - Reproducing property:

$$f \in \mathcal{F}, \; x \in \mathcal{X} \qquad f(x) = \langle k(x, \cdot), f \rangle$$

2 - Moore - Aronszajn

$$f, g \in \mathcal{F} \quad \text{RKHS}$$

$$f(x) = \sum_i \alpha_i k(x, x_i) \qquad g(x) = \sum_j \beta_j k(x, x_i)$$

w/ inner prod $\langle f, g \rangle = \sum_{ij} \alpha_i \beta_j k(x_i, x_j)$.

---

## — Basic properties and examples

**Ex (linear kernel):** $k(x, x') = \langle x, x' \rangle$ where $x, x' \in \mathbb{R}^d = \mathcal{X}$.

RKHS for $k$ :

$$\mathcal{F} = \left\{ f(x) = \sum_{i=1}^n \alpha_i k(x, x_i) : \forall n \in \mathbb{N}, \; \forall x_i \in \mathbb{R}^d, \; \forall \alpha_i \in \mathbb{R} \right\}$$

$$= \left\{ f(x) = \sum_{i=1}^n \alpha_i \langle x, x_i \rangle : \forall n \in \mathbb{N}, \; \forall x_i \in \mathbb{R}^d, \; \forall \alpha_i \in \mathbb{R} \right\}$$

$$= \left\{ f(x) = \left\langle x, \underbrace{\sum_{i=1}^n \alpha_i x_i}_{\in \mathbb{R}^d} \right\rangle : \forall n \in \mathbb{N}, \; \forall x_i \in \mathbb{R}^d, \; \forall \alpha_i \in \mathbb{R} \right\}$$

$$= \left\{ f(x) = \langle x, \theta \rangle : \theta \in \mathbb{R}^d \right\}.$$

<span style="color:purple">next lecture!</span>

Inner prod. for $\mathcal{F}$ :

$$f(x) = \langle x, \theta_1 \rangle \qquad g(x) = \langle x, \theta_2 \rangle$$

$$= 1 \cdot k(x, \theta_1) \qquad\qquad = \beta_1 k(x, \theta_2)$$

$$= \alpha_1 k(x, x_1)$$

$$\langle f, g \rangle = \langle \theta_1, \theta_2 \rangle.$$

---

**Ex (common kernels):**

1. Identity kernel : $k(x, x') = 1$. Kernel since $K_{ij} = 1$ is PSD.

2. Indicator fnc : $k(x, x') = \mathbb{1}\{\|x - x'\| \leq 0\}$. Kernel since $K = I$.

3. Polynomial kernel : $k(x, x') = (1 + \langle x, x' \rangle)^n$ is kernel.

4. Gaussian kernel : $k(x, x') = \exp\left\{ -\frac{1}{2\sigma^2} \|x - x'\|^2 \right\}$.

– Properties

1. Inner prod.: A fnc of the form $k(x, x') = \langle \phi(x), \phi(x') \rangle$

           is a kernel (see prev. lecture).

2. Summation: For $k_1$ and $k_2$ kernels, $k = k_1 + k_2$ is a kernel.

$$K_1 \succeq 0, \; K_2 \succeq 0 \; \Rightarrow \; K = K_1 + K_2 \text{ is PSD.}$$

3. Hadamard product: For $k_1$ and $k_2$ kernels, $k = k_1 \cdot k_2$ is a kernel.

$$K_1 \succeq 0, \; K_2 \succeq 0 : \qquad K_1 = \sum_k d_k u_k u_k^T \qquad K_2 = \sum_k b_k v_k v_k^T$$

$$\left( K = K_1 \circ K_2 \right)_{ij} = (K_1)_{ij} (K_2)_{ij} = \left( \sum_k d_k u_{ki} u_{kj} \right) \left( \sum_k b_k v_{ki} v_{kj} \right)$$

$$= \sum_{kl} d_k b_l \, (u_{ki} v_{li}) \cdot (u_{kj} v_{lj})$$

$$\Rightarrow \; K = \sum_{kl} d_k \cdot b_l \underset{\geq 0}{\underbrace{\phantom{xx}}} (u_k \circ v_l)(u_k \circ v_l)^T \; \succeq 0.$$

– **Polynomial kernel**:
$$k(x, x') = \underset{\text{inner prod.}}{\underbrace{(1 + \underbrace{\langle x, x' \rangle}}})^{\textcircled{m}} \to \text{product rule.}$$

<span style="color:red">sum of two kernels</span>

– **Gaussian kernel**:
$$k(x, x') = \exp\left\{ -\frac{1}{2\sigma^2} \|x - x'\|^2 \right\}$$

$$k(x, x') = \underset{k_1}{\underbrace{\exp\left\{ -\frac{\|x\|^2}{2\sigma^2} \right\} \exp\left\{ -\frac{\|x'\|^2}{2\sigma^2} \right\}}} \underset{k_2}{\underbrace{\exp\left\{ \frac{\langle x, x' \rangle}{\sigma^2} \right\}}}$$

$k_1$ is a kernel since it is an inner prod.

$k_2$ is a kernel,

$$k_2(x, x') = \exp\left\{ \frac{\langle x, x' \rangle}{\sigma^2} \right\} = \sum_{i=0}^{\infty} \frac{1}{i!} \left( \frac{\langle x, x' \rangle}{\sigma^2} \right)^i$$

$$\underset{\text{summation of poly kernels}}{\underbrace{\phantom{xxxxxxxxxxxxxxxxxxx}}}$$

$\Rightarrow \; k_1 \cdot k_2$ is a kernel.

— Learning w/ kernels

* Observe a dataset $D = \{(x_i, y_i) : i = 1 \cdots n\}$.

* We have an RKHS $\mathcal{F}$

* Consider $\quad \hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \; \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(x_i)) + \frac{\lambda}{2} \|f\|_{\mathcal{F}}^2$

**Theorem** (Representer thm): For a dataset $D = \{(x_i, y_i) : i = 1 \cdots n\}$

and a kernel $k(\cdot, \cdot)$, let $V_D = \left\{ f(x) = \sum_{i=c}^{n} x_i \, k(x, x_i) : x_i \in \mathbb{R} \right\}$.

Then, $\hat{f} \in V_D$. //

Remark: Algorithmic consequence: Minimizing over $\mathcal{F}$ = Minimizing over $V_D$.

**proof:**  — $V_D$ is a subspace of $\mathcal{F}$.

 — Define orthogonal complement of $V_D$

$$V_D^{\perp} = \left\{ f' \in \mathcal{F} : \langle f, f' \rangle = 0 \quad \forall f \in V_D \right\}.$$

(A vector space is the __sum__ of a subspace and its orthogonal complement)



$f(x) = f''(x) + f^{\perp}(x) \quad$ where

$f'' \in V_D \quad$ and $\quad f^{\perp} \in V_D^{\perp}$.
$\quad \in \mathcal{F} \qquad\qquad\qquad \in \mathcal{F}$

 — Note that $(x_i, y_i) \in D$, by the reproducing property of $\mathcal{F}$

$$f^{\perp}(x_i) = \langle \underset{\in V_D^{\perp}}{f^{\perp}}, \underset{\in V_D}{k(x_i, \cdot)} \rangle = 0.$$

$\Rightarrow f \in \mathcal{F}, \quad f(x_i) = f''(x_i) + f^{\perp}(x_i) = f''(x_i).$

$\Rightarrow \ell(y_i, f(x_i)) = \ell(y_i, f''(x_i))$

- For the regularizer: $\|f\|_{\mathcal{F}}^2 = \|f'' + f^{\perp}\|_{\mathcal{F}}^2 = \|f''\|_{\mathcal{F}}^2 + \|f^{\perp}\|_{\mathcal{F}}^2$.

- Thus, $\hat{f} = \underset{f \in \mathcal{F}}{\arg\min} \; \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(x_i)) + \frac{\lambda}{2} \|f\|_{\mathcal{F}}^2$

$$= \underset{\substack{f'' \in V_D \\ f^{\perp} \in V_D^{\perp}}}{\arg\min} \; \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f''(x_i)) + \frac{\lambda}{2} \underbrace{\|f'\|_{\mathcal{F}}^2 + \frac{\lambda}{2} \|f^{\perp}\|_{\mathcal{F}}^2}_{} \quad \geq 0.$$

<span style="color:teal">— $f^{\perp}$ has no effect.</span>

<span style="color:teal">— might as well choose $f^{\perp} = 0$.</span>

$$\Rightarrow \hat{f} \in V_D.$$

Ex (Squared error loss): For $\ell(y_i, f(x_i)) = \frac{1}{2}(y_i - f(x_i))^2$

$$\hat{f} = \underset{f}{\arg\min} \; \frac{1}{2n} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \frac{\lambda}{2} \|f\|_{\mathcal{F}}^2.$$

By the representer theorem, $\hat{f}(z) = \sum_{i=1}^{n} x_i k(x, x_i)$ for some $x_i \in \mathbb{R}$.

Finding $\hat{f}$ is equal to finding $x_i$'s.

$$\hat{x} = \underset{x \in \mathbb{R}^n}{\arg\min} \; \frac{1}{2n} \underbrace{\sum_{i=1}^{n} (y_i - \sum_{j=1}^{n} x_j k(x_i, x_j))^2}_{ii} + \frac{\lambda}{2} \underbrace{\|f\|_{\mathcal{F}}^2}_{i}$$

i) $\|f\|_{\mathcal{F}}^2 = \langle f, f \rangle = \sum_j x_i k(x_i, x_j) x_j = x^T K x$

ii) $\frac{1}{2n} \sum_{i=1}^{n} (y_i - \langle k_i, x \rangle)^2 = \frac{1}{2n} \|Y - Kx\|_2^2$

$\quad K_{ij} = k(x_i, x_j)$

$\quad K_i = \begin{bmatrix} k(x_i, x_1) \\ \vdots \\ k(x_i, x_n) \end{bmatrix}$

$$\Rightarrow \hat{x} = \underset{x \in \mathbb{R}^n}{\arg\min} \; \frac{1}{2n} \|Y - Kx\|_2^2 + \frac{\lambda}{2} x^T K x$$

$$\Rightarrow \hat{x} = \left( \frac{1}{n} K + \lambda I \right)^{-1} \frac{1}{n} Y .$$

<span style="color:purple">(More on this next lecture)</span>

**— Maximum Mean Discrepancy (MMD)**

**Goal:** Measure distance between prob. distributions given samples.

**Def:** $f, f' : X \to \mathbb{R}$, $\|f\|_\infty = \sup\limits_{x \in X} |f(x)|$ and

$$\|f - f'\|_\infty = \sup\limits_{x \in X} |f(x) - f'(x)|.$$

The following is a way to measure distance between two distributions.

**Def (MMD):** Let $p, q$ be prob. distributions on $X$. For some $\mathcal{F} = \{f : X \to \mathbb{R}\}$, define

$$d_{\mathcal{F}}(p, q) = \sup\limits_{f \in \mathcal{F}} \left| \mathbb{E}_p[f(x)] - \mathbb{E}_q[f(x)] \right|.$$

— How to choose $\mathcal{F}$? ( Want $d_{\mathcal{F}}(p, q) = 0 \Leftrightarrow p = q$ )
— How to compute $d_{\mathcal{F}}$?

**Remark:** If $\mathcal{F}$ is 1-Lipshitz fnc, $d_{\mathcal{F}}$ is $L_1$ − Wasserstein metric.

$$d_{L_1}(p, q) = W_1(p, q) \triangleq \inf \mathbb{E}\left[\|x - y\|_2\right]$$

$$\xrightarrow[\text{duality.}]{\text{By Monge-Kantorovich}}$$

couplings $(x, y)$
$x \sim p, \ y \sim q$

**Theorem (Dudley's MMD thm):** For $\mathcal{F} = C_0$ bdd cnts fncs, then $d_{C_0}(p, q) = 0 \Leftrightarrow p = q$.

— $L_1$ and $C_0$ are too complex and not practical.

**Def** (Universal kernel): A kernel $k$ is universal if its RKHS $\mathcal{F}$ is dense in $C_0$.

- $\mathcal{F}$ is dense in $C_0$ if for $f \in C_0$, $\forall \epsilon > 0$, $\exists f' \in \mathcal{F}$ s.t.
$$\|f - f'\|_\infty \le \epsilon.$$
- $\mathcal{F}$ is "representative" of $C_0$.

**Theorem** (Steinwart's thm): For unit ball $\mathcal{G} = \{f \in \mathcal{F} : \|f\|_{\mathcal{F}} \le 1\}$ where $\mathcal{F}$ is the RKHS of a universal kernel, we have
$$d_{\mathcal{G}}(p, q) = 0 \iff p = q.$$

**proof:** $\Leftarrow$ is obvious. For the other side, let $d_{\mathcal{G}}(p, q) = 0$ for some $p \neq q$.

- $p \neq q \Rightarrow d_{C_0}(p, q) = \epsilon > 0$ by Dudley's MMD thm.

$\Rightarrow \exists h \in C_0$ s.t $\left| \mathbb{E}_p[h(x)] - \mathbb{E}_q[h(x)] \right| = \epsilon.$

$\qquad\qquad\qquad\qquad\qquad\qquad$ If $f \notin \mathcal{G}$, rescale $h, f, \epsilon$ so $f \in \mathcal{G}$.

- $\mathcal{F}$ is dense in $C_0 \Rightarrow \exists f \in \mathcal{F}$ s.t $\|f - h\|_\infty \le \frac{\epsilon}{4}$

$\Rightarrow \left| \mathbb{E}_p[f(x)] - \mathbb{E}_p[h(x)] \right| \le \frac{\epsilon}{4}$ and $\left| \mathbb{E}_q[f(x)] - \mathbb{E}_q[h(x)] \right| \le \frac{\epsilon}{4}$
$\qquad\qquad * \qquad\qquad\qquad\qquad\qquad\qquad\qquad **$

- $\epsilon = \left| \mathbb{E}_p h(x) - \mathbb{E}_q h(x) \right| = \left| \mathbb{E}_p h(x) \pm \mathbb{E}_p f(x) \pm \mathbb{E}_q f(x) - \mathbb{E}_q h(x) \right|$

(by triangle ineq.) $\le \left| \mathbb{E}_p h(x) - \mathbb{E}_p f(x) \right| + \left| \mathbb{E}_p f(x) - \mathbb{E}_q f(x) \right| + \left| \mathbb{E}_q f(x) - \mathbb{E}_q h(x) \right|$
$\qquad\qquad\quad * \ \le \epsilon/4 \qquad\qquad \le d_{\mathcal{G}}(p, q) = 0 \qquad\qquad ** \ \le \epsilon/4$

$\le \frac{\epsilon}{2}$ contradiction. $\boxtimes$

- We showed that the unit ball in RKHS is good enough.
- But how to compute expectations?

* By reproducing property

$$\mathbb{E}_p f(x) = \mathbb{E}_p \langle f, k(x, \cdot) \rangle = \langle f, \mathbb{E}_p k(x, \cdot) \rangle = \langle f, \mu_p \rangle$$

$\underset{\text{fixed}}{\downarrow}$ $\underset{\text{random}}{\underbrace{\quad\quad}}$ $\quad := \mu_p$

where $\mu_p$ is the RKHS embedding of $p$.

* MMD becomes:

$$d_{\mathcal{G}}(p, q) = \sup_{f \in \mathcal{G}} |\mathbb{E}_p f(x) - \mathbb{E}_q f(x)|$$

$\mathcal{G} = \{ f : \|f\|_{\mathcal{F}} \leq 1 \}$

$$= \sup_{f \in \mathcal{G}} |\langle f, \mu_p - \mu_q \rangle|$$

$$= \|\mu_p - \mu_q\|_{\mathcal{F}} \qquad (\text{big simplification}).$$

* $$d_{\mathcal{G}}(p, q)^2 = \|\mu_p - \mu_q\|_{\mathcal{F}}^2 = \underset{1.}{\|\mu_p\|_{\mathcal{F}}^2} + \underset{2.}{\|\mu_q\|_{\mathcal{F}}^2} - \underset{3.}{2 \langle \mu_p, \mu_q \rangle}$$

1. $$\|\mu_p\|_{\mathcal{F}}^2 = \langle \mu_p, \mu_p \rangle = \langle \mathbb{E}_p k(x, \cdot), \mathbb{E}_p k(x, \cdot) \rangle$$

$$= \mathbb{E}_{p, p} [\langle k(x, \cdot), k(x', \cdot) \rangle] \qquad x, x' \sim p \text{ indep.}$$

$$= \mathbb{E}_{pp} [k(x, x')].$$

2. $\|\mu_q\|_{\mathcal{F}}^2 = \mathbb{E}_{qq} [k(y, y')]$   $y, y' \sim q$ indep.

3. $\langle \mu_p, \mu_q \rangle = \langle \mathbb{E}_p k(x, \cdot), \mathbb{E}_q k(y, \cdot) \rangle = \mathbb{E}_{pq} [k(x, y)]$

$\qquad\qquad\qquad\qquad\qquad\qquad x \sim p, y \sim q$ indep.

Plug back in:
$$d_{\mathcal{G}}(p,q)^2 = \mathbb{E}_{pp}\, k(x,x') + \mathbb{E}_{qq}\, k(y,y') - 2\mathbb{E}_{pq}\, k(x,y).$$
$$x,x' \sim p \qquad y,y' \sim q \qquad \text{indep.}$$

— Now assume $x_1 \cdots x_n \sim p$, $y_1 \cdots y_n \sim q$ indep.

**Def** ( U-statistic ):
$$U_n \triangleq \frac{1}{\binom{n}{2}} \sum_{i<j} k(x_i,x_j) + k(y_i,y_j) - k(x_i,y_j) - k(x_j,y_i)$$

**Remark** : - $U_n$ is an unbiased estimator of $d_{\mathcal{G}}(p,q)^2$.
    — It is also consistent
$$U_n \xrightarrow{P} d_{\mathcal{G}}(p,q)^2$$

    — You can use $U_n$ to measure distance between $p$ and $q$.

( Generalization: next lectures!)