# Kernel Methods Notes

- A **Kernel** is a function $K: \mathcal{X} \times \mathcal{X} \longrightarrow \mathbb{R}$ such that there exists a Hilbert space $\mathcal{H}$ and mapping $\phi: \mathcal{X} \rightarrow \mathcal{H}$ where $K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$

- A **Hilbert Space** is a vector space on which an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \longrightarrow \mathbb{R}$ is defined, this having the following properties:

  • $\langle af_1 + bf_2, g \rangle_{\mathcal{H}} = a \langle f_1, g \rangle_{\mathcal{H}} + b \langle f_2, g \rangle_{\mathcal{H}}$ (linear)

  • $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$ (symmetric)

  • $\langle f, f \rangle_{\mathcal{H}} \geq 0$, $= 0$ only when $f = 0$

- All kernels $K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$ are **positive definite** functions:

  given arbitrary $a_1, \ldots, a_n \in \mathbb{R}$, $x_1, \ldots, x_n \in \mathcal{X}$

$$\sum_i \sum_j a_i a_j K(x_i, x_j) = \sum_i \sum_j \langle a_i \phi(x_i), a_j \phi(x_j) \rangle_{\mathcal{H}}$$

$$= \left\langle \sum_i a_i \phi(x_i), \sum_j a_j \phi(x_j) \right\rangle_{\mathcal{H}}$$

$$= \left\| \sum_i a_i \phi(x_i) \right\|_{\mathcal{H}}^2 \geq 0 \quad \dashv$$

- It turns out that the opposite direction holds, as well: all positive definite functions are kernels!
- Therefore, all sums of kernels $K(x, x') = K_1(x, x') + K_2$

are kernels! for arbitrary $a_1, ..., a_n \in \mathbb{R}$, $x_1, ..., x_n \in \mathcal{X}$

$$\sum_i \sum_j a_i a_j K(x_i, x_j) = \sum_i \sum_j a_i a_j \left( K_1(x_i, x_j) + K_2(x_i, x_j) \right)$$

$$= \left\| \sum_i a_i \phi_1(x_i) \right\|^2_{\mathcal{H}_1} + \left\| \sum_i a_i \phi_2(x_i) \right\|^2_{\mathcal{H}_2} \geq$$

$$\Rightarrow \text{positive-definite} \quad \therefore \text{a kernel}$$

- All products of kernels $K(x, x') = K_1(x, x') K_2(x, x')$

are kernels!

$$K_1(x, x') K_2(x, x') = \langle \phi_1(x), \phi_1(x') \rangle_{\mathcal{H}_1} \langle \phi_2(x), \phi_2(x') \rangle_{\mathcal{H}_2}$$

can always take trace of a scalar $\Big($ $= \phi_1(x')^T \phi_1(x) \; \phi_2(x)^T \phi_2(x')$

can move a scalar into a trace $\Big($ $= \phi_1(x')^T \phi_1(x') \; \text{Trace}\left[ \phi_2(x') \phi_2(x)^T \right]$

$$= \text{Tr}\left[ \underbrace{\phi_2(x') \phi_1(x')^T}_{A^T} \; \underbrace{\phi_1(x') \phi_2(x)^T}_{B} \right]$$

Frobenius product $\Big($ $= \text{Tr}\left[ A^T B \right]$

$$= \text{vec}(A)^T \text{vec}(B)$$

$$= \left\langle \text{vec}\left( \phi_1(x') \phi_2(x')^T \right), \text{vec}\left( \phi_1(x') \phi_2(x)^T \right) \right\rangle_{\mathcal{H}}$$

$$= \langle \psi(x'), \psi(x) \rangle_{\mathcal{H}} = K(x, x') \checkmark$$

- Every kernel is associated with a unique RKHS $\mathcal{H}$, which has the following properties:
  - $\forall x \in \mathcal{X}$, $K(\cdot, x) \in \mathcal{H}$
  - $\forall x \in \mathcal{X}$, $\forall f \in \mathcal{H}$, $\langle f, K(\cdot, x) \rangle = f(x)$
  
  <u>reproducing property</u>

- Ex. RKHS defined by a <u>Fourier Series</u>

consider the space of all periodic functions on $[-\pi, \pi]$:
$$f(x) = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell e^{i\ell x}$$

We can then define the $\infty$-D space spanned by the orthonormal basis $\{e^{i\ell x}\}_{\ell=-\infty}^{\infty}$, $x \in \mathbb{R}$ together with the standard $L2$ dot product $\langle \cdot, \cdot \rangle$, to give us a Hilbert space $\mathcal{H}$, where $\langle f, g \rangle_{L2} = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \bar{\hat{g}}_\ell$.

Is $\mathcal{H}$ an RKHS? Let $K(x,y) = K(x-y)$ We check for the reproducing property: $= \sum_{\ell=-\infty}^{\infty} \hat{K}_\ell e^{i\ell(x-y)}$

$$\langle f, K(\cdot, x) \rangle_{L2} = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell \bar{\hat{K}_\ell e^{i\ell y}} \qquad = \sum_{\ell=-\infty}^{\infty} \hat{K}_\ell e^{i\ell y} e^{i\ell x}$$

$$= \sum_{\ell=-\infty}^{\infty} \hat{K}_\ell \hat{f}_\ell e^{i\ell x} \neq f(x)$$

Given this kernel, what is the dot product of the associated RKHS?

So $\mathcal{H}$ is <u>not</u> an RKHS. But we can easily modify it so that it is: $\mathcal{H}^*$ with $\langle f, g \rangle_{\mathcal{H}^*} = \frac{\sum \hat{f}_\ell \bar{\hat{g}}_\ell}{\sum_{\ell=-\infty}^{\infty} \hat{K}_\ell}$

Now, $\left\langle f, K(\cdot, x)\right\rangle_{\mathcal{H}^*} = \sum_{\ell=-\infty}^{\infty} \frac{\hat{f}_\ell \hat{K}_\ell e^{i\ell x}}{\hat{K}_\ell} = \sum_{\ell=-\infty}^{\infty} \hat{f}_\ell e^{i\ell x} = f(x)$

$\left\langle K(\cdot, x), K(\cdot, y)\right\rangle_{\mathcal{H}^*} = \sum_{\ell=-\infty}^{\infty} \frac{\hat{K}_\ell e^{-i\ell x} \hat{K}_\ell e^{i\ell y}}{\hat{K}_\ell} = \sum_{\ell=-\infty}^{\infty} \hat{K}_\ell e^{i\ell(y-x)} = K(y-x)$

Importantly, $\left\langle f, f\right\rangle_{\mathcal{H}^*} = \|f\|_{\mathcal{H}^*}^2 = \sum_{\ell=-\infty}^{\infty} \frac{|\hat{f}_\ell|^2}{\hat{K}_\ell}$, so the kernel enforces smoothness since any $f \in \mathcal{H}^*$ must have $\hat{f}_\ell$ that decay faster than $\hat{K}_\ell$ for $\|f\|_{\mathcal{H}^*}^2 < \infty$, i.e. $f(\cdot)$ must be at least as smooth (low amplitudes at higher frequencies) as $K(\cdot)$.

— **Kernel PCA** : just like normal PCA but performed in feature space, via the reproducing property:

$f^* = \underset{\substack{f \in \mathcal{H} \\ \|f\|_{\mathcal{H}}=1}}{\mathrm{argmax}}$  variance of data projected into $\mathcal{H}$ via feature map $\phi(x) = K(x, \cdot)$ along unit vector $f$

$= \underset{\|f\|_{\mathcal{H}}=1}{\mathrm{argmax}} \left\langle f, \overset{C}{\cancel{\cdots}} \right\rangle_{\mathcal{H}} = \frac{1}{N}\sum_i \left(\phi(x_i) - \frac{1}{N}\sum_j \phi(x_j)\right)^2$

$= \underset{\|f\|_{\mathcal{H}}=1}{\mathrm{argmax}} \frac{1}{N}\sum_{i=1}^{N} \left(\left\langle f, \phi(x_i)\right\rangle - \bar{\phi}\right)_{\mathcal{H}}^2 = \frac{1}{N}\sum_i^N \tilde{\phi}(x_i)^2$

$\frac{1}{N}\sum_i \left\langle f, \tilde{\phi}(x_i)\right\rangle \left\langle f, \tilde{\phi}(x_i)\right\rangle$

$\bar{\phi} = \frac{1}{N}\sum \phi(x_i)$

$\tilde{\phi}(x_i) = \phi(x_i) - \bar{\phi}$

$\frac{1}{N}\sum_i \left\langle f, \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i) f\right\rangle$

$= \underset{\|f\|_{\mathcal{H}}}{\mathrm{argmax}} \left\langle f, Cf\right\rangle, \quad C = \frac{1}{N}\sum_{i=1}^{N} \tilde{\phi}(x_i) \otimes \tilde{\phi}(x_i)$

$$\Rightarrow \frac{\partial}{\partial \delta}\left[\langle \delta, C\delta \rangle_{\mathcal{H}} + \mathcal{X}\left(\langle \delta, \delta \rangle_{\mathcal{H}} - 1\right)\right] = 0$$

$$\Longleftrightarrow C\delta = \mathcal{X}\delta$$

$$\Rightarrow \delta^* = \text{largest e-vector } C$$

$\mathcal{L}_\delta$ but this requires computing $C$, which $\not{p^A}$ lives $\cancel{\text{around}}$ in $\mathbb{R}^{\infty \times \infty}$
$\rightarrow$ How can we avoid feature space?

$\Rightarrow$ We can always express $f$ as a linear combination of data points, without loss of generality, since any dimensions orthogonal to the space spanned by $\{\phi(x_i)\}_{i=1}^n$ will disappear in the first like $\langle \delta, \phi(x_i) \rangle_{\mathcal{H}}$, thus rendering them irrelevant to the optimization:

$$f = \sum_{i=1}^{N} \alpha_i \tilde{\phi}(x_i)$$

$$\Longrightarrow f(\cdot) = \sum_{i=1}^{N} \alpha_i \tilde{K}(x_i, \cdot)$$

$\tilde{K}(x, x') = \langle \tilde{\phi}(x), \tilde{\phi}(x') \rangle_{\mathcal{H}}$

(by reproducing property)

Thus we need only solve for the $\alpha$'s:

$$Cf = \frac{1}{N} \sum_{i=1}^{N} \tilde{\phi}(x_i) \sum_{j=1}^{N} \alpha_j \langle \tilde{\phi}(x_i), \tilde{\phi}(x_j) \rangle_{\mathcal{H}}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \tilde{\phi}(x_i) \sum_{j=1}^{N} \alpha_j \tilde{K}(x_i, x_j) \Longrightarrow \langle \tilde{\phi}(x_k), Cf \rangle = \frac{1}{N} \sum_{i} \tilde{K}(x_k, x_i)$$

$\alpha_j \tilde{K}(x_k$

$$\langle \tilde{\phi}(x_k), \mathcal{X}\delta \rangle_{\mathcal{H}} = \mathcal{X} \sum_i \alpha_i \tilde{K}(x_k, x_i) \Longrightarrow \frac{1}{N} \tilde{K}\tilde{K}\alpha = \mathcal{X}\tilde{K}\alpha$$

where $\tilde{K}_{ij} = \tilde{K}(x_i, x_j)$. Since this matrix is symmetric and positive semidefinite, its inverse exists, so we get the following eigenvalue equation:

$$\tilde{K}\alpha = N\lambda\alpha$$

So we can solve for $\alpha$ by constructing the Gram matrix $\tilde{K}$ and solving the e-value equation, giving us the directions $\phi$ of greatest variance without having to work out all in feature space. (ie. biggest $\tilde{}$

Importantly, $\phi$ is a function, so kernel PCA, as opposed to regular PCA, can give us nonlinear non-linear principal subspaces rather than just hyperplanes (depending on the kernel).

— __Kernel Ridge Regression__: ridge regression in feature space

$$y = w^T\phi(x) + \epsilon, \quad \phi(x) \in \mathcal{H}$$

$$\implies w^* = \underset{w \in \mathcal{H}}{\arg\min}\left[\sum_{i=1}^{N}\left(y_i - \langle w, \phi(x_i)\rangle_{\mathcal{H}}\right)^2 + \mathcal{T}\|w\|_{\mathcal{H}}^2\right]$$

$$= \underset{w \in \mathcal{H}}{\arg\min}\left[\|Y - X^Tw\|_{\mathcal{H}}^2 + \mathcal{T}\|w\|_{\mathcal{H}}^2\right], \quad X = \begin{bmatrix}\phi(x_1) \cdots \phi(x_n)\end{bmatrix}$$

$$= \underset{w \in \mathcal{H}}{\arg\min}\left[Y^TY - 2Y^TX^Tw + w^T(XX^T + \mathcal{T}I)w\right]$$

completing the square

$$= \underset{w \in \mathcal{H}}{\arg\min}\left[Y^TY + \|(XX^T + \mathcal{T}I)^{\frac{1}{2}}w - (XX^T + \mathcal{T}I)^{-\frac{1}{2}}XY\|_{\mathcal{H}}^2 - \|(XX^T + \mathcal{T}I)^{-\frac{1}{2}}X\|\right]$$

$$= \left( X X^T + \lambda I \right)^{-1} X Y \qquad \text{\small \( \begin{pmatrix} \text{we could've done this by} \\ \text{taking derivatives, but derivatives don't} \\ \text{necessarily exist for discrete } x_i, y_i \end{pmatrix} \)}$$

To avoid having to do anything in feature space, we rewrite this in terms of the (Gram) matrix $K = X^T X$:

① Via SVD:
$$X = \begin{bmatrix} \tilde{u} \end{bmatrix} \begin{bmatrix} \tilde{S} \\ 0 \end{bmatrix} \begin{bmatrix} \tilde{V} \end{bmatrix} \qquad K_{ij} = k(x_i, x_j)$$

$\underset{D\times N}{\phantom{X}} \quad \underset{D\times D}{\phantom{[u]}} \quad \underset{D\times N}{\phantom{[S]}} \quad \underset{N\times N}{\phantom{[V]}}$

(orthogonal)  (diagonal)  (orthogonal)

Let
$$U = \tilde{u} \qquad\qquad D \times D$$
$$S = \begin{bmatrix} \tilde{S} & 0 \\ 0 & 0 \end{bmatrix} \qquad D \times D$$
$$V = \begin{bmatrix} \tilde{V} & 0 \end{bmatrix} \qquad N \times D$$

such that $X = U S V^T$

we then have:
$$w^* = \left( U S^2 U^T + \lambda I \right)^{-1} U S V^T Y$$
$$= U \left( S^2 + \lambda I \right)^{-1} U^T U S V^T Y$$
$$= U S \left( S^2 + \lambda I \right)^{-1} V^T Y$$

*can do this since $S$ is diagonal and square (hence the change from the $\tilde{}$ using SVD)*

$$= U S V^T V \left( S^2 + \lambda I \right)^{-1} V^T Y$$
$$= U S V^T \left( V^T S^2 V + \lambda I \right)^{-1} Y$$
$$= X \left( X^T X + \lambda I \right)^{-1} Y$$
$$= \underline{\underline{X \left( K + \lambda I \right)^{-1} Y}}$$

## ② Via Woodbury Identity:

$$w^* = (XX^T + \lambda I)^{-1} XY$$

$$= \left(\lambda^{-1}I - \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}X^T\lambda^{-1}\right)XY$$

$$= \left[\lambda^{-1}X - \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}\lambda^{-1}X^TX\right]Y$$

$$= \left[\lambda^{-1}X + \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}\right.$$
$$- \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}$$
$$\left. - \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}\lambda^{-1}X^TX\right]Y$$

$$= \left[\lambda^{-1}X + \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}\right.$$
$$\left. - \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}(\lambda^{-1}X^TX + I)\right]$$

$$= \lambda^{-1}X(\lambda^{-1}X^TX + I)^{-1}Y$$

$$= X\underbrace{(X^TX + \lambda I)^{-1}}_{K}Y$$

Thus, our optimal weights are a weighted sum of the data points: $w^* = \sum_i \alpha_i \phi(x_i)$, $\underline{\alpha} = (K + \lambda I)^{-1}Y$

Note that $w^*$ is a function in $\mathcal{H}$, such that its smoothness is constrained by the kernel since $\|w^*\|_{\mathcal{H}}^2 < \infty$. The larger our regularizing the constant $\lambda$, the smoother our resulting regression function $\langle w^*, \phi(x)\rangle_{\mathcal{H}} = w^*(x)$ will be.