





RKHS RFF

Hw 3 TB posted Thu: clustering

### Lecture V.1 – Build your own RKHS in 4 easy steps

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

STAT/BIOST 527 Spring 2023

### Outline

### RKHS – why bother?



Prom kernel K() to Reproducing Kernel Hilbert Space (RKHS)



Reading AoNPS Ch.: , HTF Ch.:



# Expo Random Graph Models

 $\approx x^T k x''$ • a base space X, which in ML is the input space. For example  $\mathbb{R}^d$ , or  $\{x \in \mathbb{R}^d, \|x\| \le R\}$ . • a kernel K() over X, that defines a scalar product. f(x) x ER • L<sup>2</sup>(X), the space of functions that have finite 2-norm on X.  $L^{2}(\mathsf{X}) = \{f: \mathsf{X} \to \mathbb{R}, \int_{\mathsf{X}} f^{2}(x) dx < \infty\}$ (6) A kernel defines a scalar product on X iff it is positive definite in the following sense  $\int_{\mathcal{X}} f(x)f(x')K(x,x')dxdx' > 0, \quad \text{for all } f \neq 0, f \in L^2(\mathsf{X}).$ (7)

In particular, from (7) it follows that for any set  $x^{1:n}$ , the Gram matrix

$$G = \left[ \mathcal{K}(x^{i}, x^{j}) \right]_{i,j=1}^{n} \ge 0.$$
(8)

Exercise Prove this

Marina Meila (UW Statistics)

### L V1 RKHS

From kernel K() to Reproducing Kernel Hilbert Space (RKHS)





### Remark 1

### Scalar Product

- A scalar product ( ) on X (also called inner product).
  - $\langle \ \rangle : \mathsf{X} \times \mathsf{X} \to \mathbb{R}$  is a scalar product iff it is
  - 1. Symmetric  $\langle x, x' \rangle = \langle x', x \rangle$ .
  - 2. Positive definite  $\langle x, x \rangle > 0$  for all  $x \neq 0$ .
  - 3. Bilinear (i.e. linear in each argument)  $\langle \alpha x_1 + \beta x_2, x' \rangle = \alpha \langle x_1, x' \rangle + \beta \langle x_2, x' \rangle$  (and similarly for second argument). Note that it suffices to be symmetric and linear in first argument.

 $\langle x, x' \rangle \stackrel{\text{def}}{=} x^T x' = \text{Euclidean scalar product}$ 

### The recipe

Given X, kernel K over X

• The feature map  $x \mapsto K_x() = K(x, )$ 

- Every  $x \in X$  maps to the function  $K_x : X \to \mathbb{R}$ , defined as  $K_x(u) = K(x, u)$  for all  $u \in X$ .
- Hence, each x is also a function in L<sup>2</sup>(X); we write this X → L<sup>2</sup>(X). But the set {K<sub>x</sub>, x ∈ X} has a lot of "holes", it's not useful! Must be "filled in".

**2** Start by expanding it into a linear space, the space of all finite sums of  $K_x$ 's.

$$\mathcal{H}_{0} = \text{span}\{K_{x}, x \in \mathsf{X}\} = \{\sum_{i=1}^{n} \alpha_{i} K_{x^{i}}, \text{ for } n = 1, 2, \dots, \alpha_{1:n} \in \mathbb{R}, x^{1:n} \in \mathsf{X}\}$$
(9)

This is still not enough, we would like to include limits of sequences in  $\mathcal{H}_0$ , e.g. infinite sums. For limits we need a distance.

1. map 
$$X \rightarrow function$$
  $K_X : X \rightarrow K$   
 $K_X(\underline{X'}) = K(X_1 \underline{X'})$   
 $\mathcal{A}_{-1} = \{ K_X, X \in \mathcal{X}_2 \}$  space of functions  
 $\mathcal{A} - \mathcal{H}_0 = \text{span } \mathcal{H}_{-1}$   
ring Mella (UW Statistics)  $K_X = \mathcal{K}_1 \times \mathcal{K}_2$   
 $K_X = \mathcal{K}_2 \times \mathcal{K}_3 = \mathcal{K}_1 \times \mathcal{K}_3$ 

# L V1 RKHS From kernel K() to Reproducing Kernel Hilbert Space (RKHS)



The recipe

This is still not enough, we would like to include limits of sequences in  $M_{\rm de}$  e.g. infinite same. For limits we need a distance.

### Remark 2

**Complete metric space** In a complete space  $\mathcal{H}$ , if a sequence  $\{f_n\}_{n=1}^{\infty}$  has a limit f, then f is also in  $\mathcal{H}$ ; moreover (and this is the actual definition), if a sequence is Cauchy, meaning that distance $(f_n, f_m) \to 0$  for  $m, n \to \infty$ , then the limit f exists and is in  $\mathcal{H}$ .

### Remark 3

**Hilbert space** A **Hilbert space** is an infinite dimensional vector space that has a scalar product and is complete.

# The recipe (2)

 $\langle K_x, K_{x'} \rangle_{\mathcal{H}} = K(x, x').$  Scalar frod on  $\chi$ 6 Define a scalar product  $\langle \rangle_{\mathcal{H}}$  on  $\mathcal{H}_0$ , by means of the kernel K. Let (10)

Hence, the scalar product defined by K on X, is transported to  $\mathcal{H}_0$ . This is sufficient to define the scalar product on all of  $\mathcal{H}_0$  because for any  $f, g \in \mathcal{H}_0$ ,

$$f, g\rangle_{\mathcal{H}} = \langle \sum_{i=1}^{n} \alpha_{i} K_{u^{i}}, \sum_{j=1}^{m} \beta_{j} K_{v^{j}} \rangle_{\mathcal{H}} = \sum_{i=1}^{n} \sum_{j=1}^{m} \alpha_{i} \beta_{j} \langle K_{u^{i}}, K_{v^{j}} \rangle_{\mathcal{H}}$$
(11)  
$$= \sum_{i=1}^{n} \sum_{j=1}^{m} \alpha_{i} \beta_{j} K(u^{i}, v^{j})$$
(12)

Exercise Prove that  $\langle \rangle_{\mathcal{H}}$  is a scalar product. o define a norm  $\|f\|_{\mathcal{H}}^{2} = \langle f, f \rangle_{\mathcal{H}} \stackrel{\text{distance}}{\longrightarrow} \left( \begin{array}{c} f, g \end{array} \right) = \\ \|f\|_{\mathcal{H}}^{2} = \langle f, f \rangle_{\mathcal{H}} \stackrel{\text{distance}}{\longrightarrow} \left( \begin{array}{c} f \\ g \end{array} \right) = \\ \|f - g\|_{\mathcal{H}} \\ \|f - g\|_{\mathcal$ The scalar product  $\langle \rangle_{\mathcal{H}}$  allows us to define a norm

Now we can complete  $\mathcal{H}_0$  to  $\mathcal{H}$ .

Voila!  $\mathcal{H}$  is your **Reproducing Kernel Hilbert Space (RKHS)**.

(13)

### Recipe, summarized

Input X, kernel K Map X  $\hookrightarrow L^2(X)$  by the feature map  $x \mapsto K_x() = K(x, )$ Make it a linear space  $\mathcal{H}_0 = \operatorname{span}\{K_x, x \in X\} = \{\sum_{i=1}^n \alpha_i K_{x^i}, \text{ for } n = 1, 2, \dots, \alpha_{1:n} \in \mathbb{R}, x^{1:n} \in X\}$ Define scalar product  $\langle \rangle_{\mathcal{H}}$  on  $\mathcal{H}_0$ , by  $\langle K_x, K_{x'} \rangle_{\mathcal{H}} = K(x, x')$ . Complete  $\mathcal{H}_0$  to  $\mathcal{H}$  using  $|| \, ||_{\mathcal{H}}$ . K on  $X \times X$ reproduced on  $\mathcal{H}X\mathcal{H}$ 

# The name RKHS explained

- Reproducing Kernel Hilbert Space
  - means the space of functions has a scalar product and is complete
- Reproducing Kernel Hilbert Space
  - the scalar product comes from a kernel
- Reproducing Kernel Hilbert Space
  - in addition, this space has the Reproducing property (coming next!)

#### Properties of RKHS's

# The Reproducing Property $\langle f, K_x \rangle_{\mathcal{H}} = f(x)$

• Let's prove it. Remember  $f(x) = \sum_{i=1}^{n} a_i K_{u_i}(x)$  for  $f \in \mathcal{H}_0$ ,  $x \in X$ .

$$(f, \underline{K_x})_{\mathcal{H}} = \sum_{i=1}^n a_i \langle K_{u_i}, K_x \rangle_{\mathcal{H}}$$
 (16)

$$f(\underline{x}) = \sum_{i=1}^{n} a_i K(u_i, x) = f(x)$$

$$= \sum_{i=1}^{n} a_i K(u_i, x) = f(x)$$
(17)

- In other words, if we map x into  $\mathcal{H}$  by  $x \mapsto K_x$  and calculate the scalar product with some  $f \in \mathcal{H}$ , the result is the same as applying f to  $x \in X$ .
- One can say that  $K_x$  reproduces x
- Or alternatively that *f* ∈ *H*, by Riesz's Theorem, defines the linear functional ⟨*f*, ⟩<sub>*H*</sub>. This functional on *H* reproduces the effect of *f* on X.

# Mercer's Theorem

# Canonical basis for L2(X)

• Define the transport operator  $T: L^2(X) \to \mathcal{H}$ 

$$Tf = \int_{X} f(u)K(u,u) du \quad \Leftrightarrow \quad Tf(x) = \int_{X} f(u)K(u,x) du$$
(18)

( like convolution )

• Let  $\{(\lambda_i, \psi_i)\}_i$  be the eigenvalue, eigenfunction pairs of T • The Mercer Theorem says that, under certain conditions on X and K, the operator T

- has a discrete spectrum,
- 2) is positive semidefinite  $\lambda_i > 0$  for i = 1, 2, ...
- **(**) the eigenfunctions  $\{\psi_i\}_{i=1}^{\infty}$  form an orthogonal basis for  $L^2(X)$  $\begin{array}{c} x \xrightarrow{\mu} K_{x} \equiv K(x, \cdot) \\ f \xrightarrow{\mu} Tf \quad (Tf) \in \widehat{f} \quad \int f(u) \quad K(x, u) \, du = \widehat{f}(x) \\ f \xrightarrow{\mu} f(x) \quad f \xrightarrow{\mu}$

H

H

### L V1 RKHS Properties of RKHS's

2023-05-08

Mercer's Theorem

#### Mercer's Theorem



Ans a decrete spectrum,
 Is positive semidefinite 3i ≥ 0 for i = 1, 2, ...
 the restriction semidefinite 3i ≥ 0 for i = 1, 2, ...
 the restriction semidefinite 3i ≥ 0 for m an orthogonal basis for L<sup>2</sup>(X)

#### Remark 4

#### (Linear) Operator

- An (linear) operator T is a (linear) function from a space of functions to another.
- For example the derivative maps a function  $f : \mathbb{R} \to \mathbb{R}$  to its derivative f'; we can write that derivative :  $\mathcal{C}^1(\mathbb{R}) \to \mathcal{C}^0(\mathbb{R})$  is a linear operator.
- For an operator T and function f, we denote by  $g = Tf \equiv T(f)$ , the function resulting from applying T to f.
- Furthermore, if we calculate this function g at point x, we write g(x) = Tf(x)
- Operators have eigenfunctions and eigenvalues defined as  $T\psi = \lambda\psi$  for some  $\lambda \in \mathbb{R}$
- The set of eigenvalues  $\{\lambda, \text{ such that } T\psi = \lambda\psi \text{ for some }\psi\}$  is the spectrum of T.
- The spectrum of an operator is usually more complicated than the spectrum of a matrix; for example, it can contain continuous intervals, the whole real line, limit points. If the spectrum contains none of these, i.e. consists of only isolated eigenvalues, we say the spectrum is discrete.

# The feature map revisited

• The consequences of this theorem are remarkable. In particular, it lets us express the kernel itself in the basis of T.

$$\mathcal{K}(\underline{x},\underline{x}') = \sum_{i=1}^{\infty} \lambda_i \underline{\psi_i(\underline{x})} \underline{\psi_i(\underline{x}')}$$
(19)

• Therefore,

$$K(x,x) = \sum_{i=1}^{\infty} \lambda_i \psi_i(x)^2.$$
 (20)

• From here, it is easy to see that the feature map  $x \mapsto K_x$  can also be written as

$$x \mapsto \left[\sqrt{\lambda_i}\psi_i(x)\right]_{i=1}^{\infty} \equiv \text{feature map in} (21)$$

$$\begin{bmatrix} \psi_i \end{bmatrix} \text{ basis}$$

• And finally, the infinite sum converges uniformly

$$\lim_{\chi \to \infty} \sup_{x,x'} \left| K(x,x') - \sum_{i=1}^{\mathcal{M}} \lambda_i \psi_i(x) \psi_i(x') \right| = 0$$
(22)





- It's important to remember that there are 2 scalar products here. There is the scalar product induced by the kernel K on H, defined in (10) and (11), and there is the standard scalar product on L<sup>2</sup>(X) defined by ⟨f,g⟩ = ∫<sub>X</sub> f(x)g(x)dx.
- The basis  $\{\psi_i\}$  is orthonormal w.r.t. the  $L^2(X)$  scalar product.

How to prove (19).

•  $\psi_j$  is eigenfunction, hence

$$\int_{\mathsf{X}} \psi_j(\mathbf{x}') \mathcal{K}(\mathbf{x}, \mathbf{x}') d\mathbf{x}' = \lambda_j \psi_j(\mathbf{x}).$$
<sup>(23)</sup>

- Now  $K_x$  itself has a decomposition in the basis,  $K_x = \sum_i \gamma_i(x)\psi_i$ , where  $\gamma_i(x)$  are the coefficients.
- Let's plug this decomposition in (23)

$$\int_{\mathsf{X}} \psi_j(\mathbf{x}') \mathcal{K}(\mathbf{x}, \mathbf{x}') d\mathbf{x}' = \int_{\mathsf{X}} \psi_j(\mathbf{x}') \sum_i \gamma_i(\mathbf{x}) \psi_i(\mathbf{x}') d\mathbf{x}'$$
(24)

$$= \sum_{i} \gamma_{i}(x) \int_{\mathsf{X}} \psi_{j}(x') \psi_{i}(x') dx'$$
 (25)

$$= \gamma_j(x) = \lambda_j \psi_j(x).$$
 (26)

• Hence  $K_x(x') \equiv K(x, x') = \sum_i \lambda_i \psi_i(x) \psi(x')$ . Done.

Kurnel SVM  
Given 
$$\mathfrak{D} = f(x_{i}^{i}, y_{i}^{i}), i = 1:n \mathfrak{Z}, C = regularization param.
Kurnel  $K(r, 1 > 0$   
1. Compute  $G = [K(x_{i}^{i}, x_{i}^{i})]_{ij=1:n}, G = [y_{i}^{i}y_{i}^{i}K(x_{i}^{i}, x_{i}^{i})]_{ij=1:n}$   
2. Solve dual SVM  
 $Max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $max \quad 1^{T}\alpha - \frac{1}{4}\alpha^{T}G\alpha \quad st. \quad \alpha \in [0, C], \quad \sum_{i=1}^{n} y_{i}^{i}\alpha_{i}=0$   
 $dim = n$   
 $dim = n$   
 $dim = n$   
 $f(x) = \sum_{i=1}^{n} \alpha_{i}y_{i}K(x_{i}, x_{i}) + b$   
 $Representer theorem$   
 $f(x) = \sum_{i=1}^{n} \alpha_{i}y_{i}K(x_{i})$   
 $(a) \quad f represented by \quad data''$   
 $im \quad H : \quad f = \sum_{i=1}^{n} \alpha_{i}y_{i}Kx_{i}$   
 $Kx_{i}$$$

Kernel Machines in practice X not enclidean  $\leftarrow$  X(,) Strings: DNA sequences profeirs phylogenetic trees language kernel - machines.org

GROKT : Can now use SVM, kunel repression, ... on strings, trees, ...

SVM for hig data

### Lecture VI-2: SVM with Random Fourier Features

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

Spring, 2023

Reading: Ali Rahimi and Ben Recht "Random features for large-scale Kernel Machine", NIPS 2007. Test of Time Award, NIPS 2017.

# Problem: Kernel machines scale with sample size n

- ► Gram matrix  $G = [k(x^i, x^j)]_{i,j=1}^n \in \mathbb{R}^{n \times n}$  Expensive/intractable for *n* large!
- Want to: benefit from infinite dimensional feature spaces, e.g. Gaussian kernel, AND have constant dimension <u>D</u> for any n
- Idea approximate  $\overline{k(x, x')}$  with finite sum.
- Equivalently, approximate feature space H with D-dimensional feature space. How? Pick D features at random!

primal SVM

min  $\frac{1}{4} \|W\|^2 + C \geq 3i$   $W_{vb}i_{50} = s.t yi[WTXi+b] \geq (-3)$  $P^{D} R R^{h} = 3i \geq 0$ 

### Why is this possible? Bochner's Theorem



Let K(x, x') = K(x - x') be a continuous shift invariant kernel.

### Theorem [Bochner]

 $\frac{K(x - x')}{\max p(\omega)}$  is a positive definite kernel iff K(z) is the Fourier transform of some non-negative measure  $p(\omega)$ .  $\frac{K(z) = \int p(\omega)e^{-i\omega^{T}z}d\omega \qquad k = Fourier(P) \qquad (1)$ 

$$e^{-i\omega^{T}(x-x')} = K(x-x') = \int_{\mathbb{R}^{d}} e^{|\omega|^{2}/2} d\omega = \int_{\mathbb{R}^{d}} e^{|\omega|^{2}/2} d\omega$$

# From Bochner to RFF

- Note that  $e^{-i\omega z} = e^{-i\omega^T x} (e^{-i\omega^T x'})^*$  and let  $\zeta_{\omega}(x) = e^{-i\omega^T x}$
- ► Then  $K(z) = E_{p(\omega)}[\zeta_{\omega}(x)\zeta_{\omega}^*(x')] \approx \frac{1}{D} \sum_{i=1}^{D} \zeta_{\omega_i}(x)\zeta_{\omega_i}^*(x')$  with  $\omega_{1:D} \sim \text{i.i.d. } p(\omega)$
- D is the sample size, must be large enough for good approximation
- $\blacktriangleright \zeta_{\omega_{1},D}$  form a random feature space of dimension D
- Feature map is  $x \to \tilde{\phi}(x) = \frac{1}{\sqrt{D}} [\zeta_{\omega_1} \dots \zeta_{\omega_D}]$

**Fact** Because K() is real, the random complex features  $\zeta_{\omega} \leftarrow \sqrt{2}cos(\omega^T x + \omega_0)$  with  $\omega_0 \sim uniform[0, 2\pi]$ 

- **Significance** Infinite dimensional feature vector  $\phi(x)$  approximated by D-dimensional feature vector  $\tilde{\phi}(x)$ . Hence, primal problem of dimension D can be solved instead of dual of dimension *n*.
- Opens up SVM/kernel machines for large data

# Approximation

### Theorem [Rahimi and Recht 07]

Assume space  $\mathcal{X}$  is compact of diameter  $d_{\mathcal{X}}$  and let  $\sigma_p^2 = E_p[\omega^T \omega]$  be the standard deviation of  $p(\omega)$ . Then,

$$Pr\left[\sup_{x,x'\in\mathcal{X}}|\tilde{\phi}(x)^{T}\tilde{\phi}(x')-\mathcal{K}(x,x')|\geq\epsilon\right]\leq e^{-\frac{D\epsilon^{2}}{4(d+2)}}\left(\frac{2^{4}\sigma_{P}d_{\mathcal{X}}}{\epsilon}\right)^{2}$$
(2)

2. For  $\delta$  confidence level,

$$D = \Omega\left(\frac{d}{\epsilon^2}\ln\frac{\sigma_p d_{\mathcal{X}}}{\epsilon}\right)$$
(3)

1.

Spring, 2023

# Spring, 2023

### Kernel machine with RFF algorithm

In Data  $x^{1:n}, y^{1:n}$ , kernel K

- 1. Fourier transform  $p(\omega) = \frac{1}{2\pi} \int_{\mathbb{R}^d} e^{-i\omega^T z} K(z) dz$ .
- 2. Choose D.
- 3. Sample  $\omega_{1:D}$  i.i.d. from *p*. Sample  $\omega_{0,1:D}$  uniformly from  $[0, 2\pi]$ .
- 4. Map data to features  $\tilde{\phi}(x^i) = \sqrt{\frac{2}{D}} [\cos(\omega_i^T x^i + \omega_{0,j})]_{j=1:D}$  for all i = 1: n.
- 5. Solve SVM Primal problem; obtain  $w \in \mathbb{R}^D$  and intercept  $b \in \mathbb{R}$ .