

# ASSIGNMENT 1 - MVP's

LINDI LI  
3460570

GEWEI CAO  
3461232

December 2022

## § 1 INTRUDCTION

blabla

## § 2 REPRODUCING KERNEL HILBERT SPACE

### § 2.1 HILBERT SPACE

Recall that an inner product  $\langle a, b \rangle$  can be

- a usual dot product:  $\langle a, b \rangle = a'b = \sum_i a_i b_i$ .
- a kernel product:  $\langle a, b \rangle = k(a, b) = \psi(a)' \psi(b)$ , where  $\psi(a)$  may have infinite dimensions.

A normed space is a vector space  $N$  on which a norm is defined. A nonnegative function  $\|\cdot\|$  is a norm if and only if  $\forall f, g \in N$  and  $\alpha \in \mathbb{R}$ :

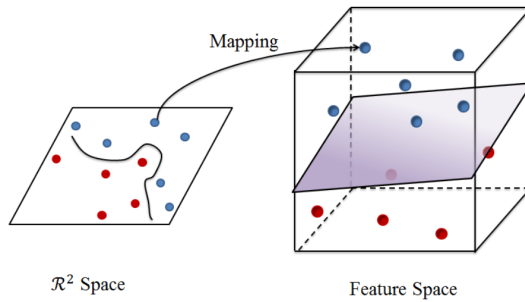
- $\|f\| \geq 0$  and  $\|f\| = 0$  iff  $f = 0$ ;
- $\|f + g\| \leq \|f\| + \|g\|$ ;
- $\|\alpha f\| = |\alpha| \|f\|$

Now we can define a Hilbert space an inner product space that is complete and separable with respect to the norm defined by inner product. An example of a defined norm in Hilbert space (i.e, the space  $L_2$  of square integrable functions) can be

$$\|f\| = \left( \int_a^b f^2(t) dx \right)^{\frac{1}{2}}. \quad (1)$$

### § 2.2 KERNEL METHOD

The motivation of kernel method is simple. Imagine there are some blue dots and red dots on a vector space  $\mathbb{R}^n$  and we want to separate them by color. As it shows in the left hand side figure, it is difficult to divide them through a straight line. However, we may be able to separate them easily by mapping each dot into a high-dimension feature space. The figure below shows how mapping works:



Let's now use a simple example to illustrate the idea of kernel method. we set two vectors  $\mathbf{x} = [x_1 \ x_2]$  and  $\mathbf{y} = [y_1 \ y_2]$  in a two-dimension space. Two functions  $\phi(x)$  and  $\phi(y)$  are defined as:

$$\phi(x) = [x_1 x_1 \quad x_1 x_2 \quad x_2 x_1 \quad x_2 x_2]$$

$$\phi(y) = \begin{bmatrix} y_1 y_1 & y_1 y_2 & y_2 y_1 & y_2 y_2 \end{bmatrix}$$

We now successfully mapping them into a four-dimension feature space. Therefore, instead of computing the inner product of  $\langle \phi(x), \phi(y) \rangle$ , we can define a corresponding kernel function  $K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle^2$ . It can be easily proofed that  $K(\mathbf{x}, \mathbf{y})$  is the same as  $\langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$ :

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{x}, \mathbf{y} \rangle^2 \\ &= (x_1 y_1 + x_2 y_2)^2 \\ &= x_1^2 y_1^2 + 2x_1 y_2 x_2 y_1 + x_2^2 y_2^2 \\ &= \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle \end{aligned}$$

To write the above example in a general form, we first set an equation where  $\phi(\cdot) \in \mathbb{R}^m$  and  $\phi(x)$  is defined as the mapping function:

$$\begin{aligned} y &= \phi(x)^\top w \\ &= \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_m(x) \end{bmatrix}^\top w \\ &= [\phi_1(x) \cdots \phi_m(x)] w \end{aligned} \tag{2}$$

$Y$  and  $\Phi$  in generalization is defined by:

$$Y = \begin{bmatrix} y_1 & \cdots & y_n \end{bmatrix}^\top \tag{3}$$

$$\begin{aligned} \Phi &= [\phi(x_1) \quad \cdots \quad \phi(x_n)]^\top \\ &= \begin{bmatrix} \phi_1(x_1) & \cdots & \phi_m(x_1) \\ \vdots & \vdots & \vdots \\ \phi_1(x_n) & \cdots & \phi_m(x_n) \end{bmatrix} \end{aligned} \tag{4}$$

In most of the cases, we want to avoid computing  $\phi(x)$  in a explicit way, especially when  $m$  is very large. Therefore, we define a kernel function:

$$[\Phi \Phi^{top}]_{i,j} = \phi(x_i)^{top} \phi(x_j) = K(x_i, x_j) \tag{5}$$

$\phi(x_i)^{top} \phi(x_j)$  is usually very difficult to compute, and thus we use  $\mathbb{K}(x_i, x_j)$  to replace. In the simplest form of machine learning, in order to predict  $x$ , the algorithm collects the samples in the training set  $\chi$  that are similar to  $x$ , and then take the weighted value of these samples as the predict value of  $x$ . Here comes the questions:

- How to measure the similarity between samples?
- How to weight the value of each sample?

In general, the higher the similarity of the sample to our point of interest  $x$ , the more the sampling weights. To evaluate the similarity between two observations, a kernel is defined as a function of two input patterns  $k(x_i, x_j)$ , mapping onto a real-valued output. For example, the Gaussian kernel is

$$k(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{\sigma^2}}, \quad (6)$$

where  $\|x_i - x_j\|$  is the Euclidean distance between  $x_i$  and  $x_j$ , and  $\sigma^2 \in \mathbb{R}^+$  is the bandwidth of the kernel function.

We now define that  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel if

- $k$  is symmetric:  $k(x, y) = k(y, x)$ .
- $k$  is positive semi-definite, meaning that  $\sum_i \sum_j \alpha_i \alpha_j k(x_i, x_j) \geq 0, \forall \alpha_i, \alpha_j \in \mathbb{R}, x \in \mathbb{R}^D, D \in \mathbb{Z}^+$ .
- We define the corresponding kernel matrix as the matrix  $K$  with entries  $k_{ij} = k(x_i, x_j)$ , the second property of  $k$  is equivalent to saying that  $c'Kc \geq 0$ .

From the similarity-based point of view, the use of kernels for regression can be described in two stages. We set  $y_i \in \mathbb{R}$  as dependent variable, and  $x_i$  as a  $1 \times D$  vector  $x_i \in \mathbb{R}^D$  of covariate values. Assume that  $(y_i, x_i)$  where  $i = 1, \dots, N$  is i.i.d. We first define a target function  $y = f(x)$  and assume that in a space of functions, there exists a function that can estimate  $y = f(x)$  well. The target function could be represented by

$$f(x) = \sum_{i=1}^N w_i k(x, x_i), \quad (7)$$

A perfect solution to estimate  $N$  parameters for  $N$  observations would be choosing  $\hat{w} = K^{-1}y$ . However, such a fit could result in extremely high variance. Therefore, we impose an additional assumption that a smoother curve with less oscillations is preferred. We then utilize regularization to simplify the function and satisfy the additional assumption in the second stage. To achieve this purpose, Hilbert space and reproducing kernel Hilbert space will be introduced below.

### § 3 BAYIES

blabla

## § 4 EMPIRICAL STUDY

blabla