# Assignment 1 - MVP's

Lindi Li            Gewei Cao
3460570             3461232

December 2022

# § 1 Intrudction

blabla

# § 2  Reproducing Kernel Hilbert Space

## § 2.1  What is Kernel

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. We set $y_i \in \mathbb{R}$ as dependent variable, and $x_i$ as a $1 \times D$ vector $x_i$ in $\mathbb{R}^D$. Assume that $(y_i, x_i)$ where $i = 1, \ldots, N$ is i.i.d. To evaluate the similarity between two observations, a kernel is defined as a function of two input patterns $k(x_i, y_i)$, mapping onto a real-valued output. For example, the Gaussian kernel is

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

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 $k : \chi \times \chi \to \mathbb{R}$ is a kernel if

- $k$ is symmetric: $k(x, y) = k(y, x)$.

- $k$ is positive semi-definite, meaning that $\sum_i \sum_j \alpha_i \alpha_j k(x_i, x_j) \geq 0, \forall \alpha_i, \alpha_j \in \mathbb{R}, x \in \mathbb{R}^{\mathbb{D}}, D \in \mathbb{Z}^+$.

From the similarity-based point of view, the use of kernels for regression can be described in two stages. We first set 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 is represented by

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

In the second stage, we utilize regularization to simplify the function. To achieve this purpose, Hilbert space and reproducing kernel Hilbert space will be introduced below.

## § 2.2  Hilbert Space

Recall that an inner product $< a, b >$ can be

- a usual dot product: $< a, b >= a'b = \sum_i a_i b_i$.

- a kernel product: $< a, b >= k(a, b) = \psi(a)' \psi(b)$, where $\psi(a)$ may have infinite dimensions.

We define a Hilbert space an inner product space that....

# § 3 GAUSSIAN PROCESS AND BAYESIAN PERSPECTIVE

## § 3.1 DEFINITION

We have data $\mathbf{D} = \{(x_i, y_i)\}_{i=1}^M$, and assume that mean of y is 0.

Task: find the distribution of $f^*(x)$.

Assume that the true form of prediction function is: $y_i = f(x_i) + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$. Here we have a M dimensional dependent variable $\mathbf{y}$, and a $M \times N$ dimensional independent variable $\mathbf{X}$, where M is the number of observations, and N is the dimension of x, i.e. $x_i \in \mathbb{R}^N$. The function $f(x_i) : \mathbb{R}^N \to \mathbb{R}$ takes vector $x_i \in \mathbb{R}^N$. Let $K_{X,X} = k(x, x^T)$ which is the matrix of $k(x_i, x_j)$. Thus, $\mathbf{K}$ is a $M \times M$ matrix.

The assumption of Gaussian Process is as following:

For a given vector $\mathbf{y}$, and its corresponding data $\mathbf{X}$, where vector $y \in \mathbb{R}^M$ and $X$ is $M \times N$ matrix. In addition, for $\mathbf{y}$ and $\mathbf{X}$ data, the error term $\epsilon \sim \mathcal{N}(0, \Sigma^\epsilon)$, and $\Sigma^\epsilon = diag(\sigma_1^2, \sigma_2^2, \sigma_3^2, ......, \sigma_M^2)$. Meanwhile we have arbitrary $n \times N$ matrix $\mathbf{Z}$ and predicted value $f^*(z) \in \mathbb{R}^n$, where $z = (z_1, z_2, z_3, ......, z_n)^T$.

Then we assume $\mathbf{y}$ and $f^*(z)$ follow a $(M + n)$ multivariate normal distribution(MVN):

$$\begin{bmatrix} f^*(z) \\ y \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \mu_{f^*(z)} \\ \mu_y \end{bmatrix} , \begin{bmatrix} K_{Z,Z} & K_{Z,X} \\ K_{X,Z} & \hat{K}_{X,X} \end{bmatrix} \right) \tag{1}$$

where $\hat{K}_{X,X} = K_{X,X} + \Sigma^\epsilon$.

Then given data $\mathbf{y}$, $\mathbf{X}$ and $\mathbf{Z}$, according to the conditional distributions of the multivariate normal distribution[1], we have the posterior distribution

$$f^*(z) | y, X, Z \sim \mathcal{N}(\mu_{f^*(z)} + K_{Z,X} \hat{K}_{X,X}^{-1}(y - \mu_y), K_{Z,Z} - K_{Z,X} \hat{K}_{X,X}^{-1} K_{X,Z}) \tag{2}$$

## § 3.2 INTUITION BEHIND GAUSSIAN PROCESS

The idea behind this process is that, assume our interested function is $f(x)$, $f(x) : \mathbb{R}^N \to \mathbb{R}$, and we have an arbitrary vector of independent variable $x = (x_1, x_2, ......, x_M)^T$, and for each $x_i, i = 1, 2, ..., M, x_i \in \mathbb{R}^N$, then we can obtain a series of $f(x) = (f(x_1), f(x_2), f(x_3), ......, f(x_M))^T$. We assume that the series of f(x) follows a multivariate normal distribution which is:

$$f(x) \sim \mathcal{N}(\mu(x), k(x, x^T)) \tag{3}$$

This is the prior distribution of our function $f(x)$, here we have a set of infinitely functions that follow this distribution, their mean is the function $\mu(x_i)$, and the variance of them is $k(x_i, x_i^T)$. This makes the distribution of $f(x)$ to be called Gaussian Process (GP). Note that if we add a noise term $\epsilon \sim \mathcal{N}(0, \Sigma^\epsilon)$, then our prior distribution of

---

[1]https://statproofbook.github.io/P/mvn-cond

$y = f(x) + \epsilon \sim \mathcal{N}(\mu(x), k(x, x^T) + \Sigma^\epsilon)$ is also a Gaussian Process. Here we use kernel matrix to denote variance-covariance matrix because kernel value represents how near two data points in the space are, with this property we can obtain a smooth function.

Remind that our goal is to estimate the distribution of $f(x^*)$ given observed training data set $D = \{x_i, y_i\}_{i=1}^M$ and test data set $\{x_j^*\}_{j=1}^n$. Firstly we compare our nonparametric case to a parametric case. In a parametric case, assume the parameter $\theta$ determines the form of $f_\theta(\cdot)$, according to the Bayesian rule, $p(y^*|x^*, x, y) = \int_\theta p(y^*, \theta|x^*, x, y)d\theta = \int_\theta p(y^*|\theta, x^*)p(\theta|x, y)d\theta$, where $y^*$ is the prediction of given data $x^*$, and its form of model is determmined by paramater $\theta$. Estimated $\theta$ value is determined by training data $D$. This is to say that we update our parameter $\theta$ by given $D$, and use $p(\theta|x, y)$ as a new prior probability, and based on this to predict posterior of $y^*$.

Therefore, back to our GP nonparametric case, $\theta$ could be substituted by function $f(\cdot)$. One can show that the joint distribution of $(f(x^*), y)^T$ follows a multivariate normal distribution as in the definition before, because of the assumption of GP and the property of MVN. With the joint distribution, we want to find posterior probability: $p(f(x^*)|x^*, x, y) = \int p(f(x^*)|f, x^*)p(f|x, y)df$, where $p(f|x, y)$ is the posterior of $f(\cdot)$ given $D$, and is regarded as prior when estimating $p(f(x^*)|x^*, D)$, this process is called Bayesian updating. Fortunately, we do not need to take any integral in GP, becaus the posterior of $f(x^*)$ could be calculated by formula of conditional distributioin in MVN as mentioned in former section.

## § 3.3 Gaussian Process in research paper

In our object paper, for given price data P,

## § 4  Empirical study

blabla