

Random Fourier Features for Gaussian Process Model

Tetsuya Ishikawa

tiskw111@gmail.com

Abstract

This article describes the procedure for applying random Fourier features [1] to Gaussian process model [2]. This makes it possible to speed up the training and inference of Gaussian process model, and to apply the model to large-scale data.

Gaussian process model [2] is one of the supervised machine learning frameworks designed on a probability space, and is widely used for regression and classification tasks, as well as support vector machine and random forest. The major difference between Gaussian process model and other machine learning models is that Gaussian process model is a *stochastic* model. In other words, since Gaussian process model is formulated as a stochastic model, it can provide not only the predicted value but also a measure of uncertainty for the prediction. This is a very useful property that can improve the explainability of machine learning model.

On the other hand, Gaussian process model is also known for its high computational cost of training and inference. If the total number of training data is $N \in \mathbb{Z}^+$, the computational cost required for training is $O(N^3)$, and the computational cost required for inference is $O(N^2)$, where O is *Bachmann–Landau notation*. The problem is that the computational cost is given by a power of the total number of training data N , which can be an obstacle when applying the model to large-scale data. This comes from the fact that Gaussian process model has the same mathematical structure as kernel methods, in other words, the kernel support vector machine also has the same problem.

One of the methods to speed up kernel methods is random Fourier features [1] (hereinafter abbreviated as RFF). This method can significantly reduce the computational cost while keeping the flexibility of kernel methods by approximating a kernel function as an inner product of finite dimensional vectors. Specifically, the computational cost required for training can be reduced to $O(ND^2)$, and the amount of calculation required for inference can be reduced to $O(D^2)$, where $D \in \mathbb{Z}^+$ is a hyperparameter of RFF and can be specified independently of the total number of training data N . Since Gaussian process model has the same mathematical structure as kernel methods, RFF can be applied to Gaussian process model as well. This evolves Gaussian process model into a more powerful, easy-to-use, and highly reliable ML tool.

However, when applying RFF to Gaussian process model, some mathematical techniques are required that are not straightforward. Unfortunately, there seems to be no article in the world that mentions its difficulties and solutions, so I decided to leave this document. If you prefer the Japanese version of this document, see the repository [3].

1 Gaussian Process Model Revisited

This section gives an overview of Gaussian process model. Unfortunately, this section can not cover details such as the for-

mulation and derivation of Gaussian process models due to the limitation of pulp, so if you are interested in the details, please refer [2].

Let $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ be training data, and $\sigma \in \mathbb{R}^+$ be a standard deviation of the label observation error, where $\mathbf{x}_n \in \mathbb{R}^M$, $y_n \in \mathbb{R}$. Gaussian process model describes the prediction as a probability variable that follows the normal distribution. If the test date is $\xi \in \mathbb{R}^M$, the expectation of the prediction is given by:

$$m(\xi) = \widehat{m}(\xi) + (\mathbf{y} - \widehat{m})^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\xi), \quad (1)$$

and the covariance of the test data $\xi_1, \xi_2 \in \mathbb{R}^M$ is given by:

$$v(\xi_1, \xi_2) = k(\xi_1, \xi_2) - \mathbf{k}(\xi_1)^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\xi_2), \quad (2)$$

where the function $k : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$ is a kernel function, the matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$ is a kernel matrix defined as

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}, \quad (3)$$

and the vector function $\mathbf{k}(\xi) : \mathbb{R}^M \rightarrow \mathbb{R}^N$ and the vector $\mathbf{y} \in \mathbb{R}^N$ is defined as

$$\mathbf{k}(\xi) = \begin{pmatrix} k(\xi, \mathbf{x}_1) \\ \vdots \\ k(\xi, \mathbf{x}_N) \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}, \quad (4)$$

respectively. Also, $\widehat{m}(\xi)$ is the prior distribution of the prediction, and $\widehat{m} = (\widehat{m}(\mathbf{x}_1), \dots, \widehat{m}(\mathbf{x}_N))^\top$ is the prior distribution of the predicted values of the training data. If you don't need to set prior distribution, it's common to set $\widehat{m}(\cdot) = 0$ and $\widehat{m} = \mathbf{0}$.

You can compute the variance of the prediction of the test data ξ by substituting $\xi_1 = \xi_2 = \xi$ into the equation (2),

$$v(\xi, \xi) = k(\xi, \xi) - \mathbf{k}(\xi)^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\xi). \quad (5)$$

2 RFF Revisited

In this section, we revisit random Fourier features. Unfortunately, this article don't have enough space to explain the details as the same as the previous section, therefore if you would like to know more details, please refer to the original paper [1].

Let the function $k : \mathbb{R}^M \times \mathbb{R}^M \rightarrow \mathbb{R}$ be a kernel function. In RFF, the kernel function can be approximated as

$$k(\mathbf{x}_1, \mathbf{x}_2) \simeq \boldsymbol{\phi}(\mathbf{x}_1)^\top \boldsymbol{\phi}(\mathbf{x}_2), \quad (6)$$

where $\boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^D$ is a feature vector extracted from the data \mathbf{x} and $D \in \mathbb{Z}^+$ is the dimension of $\boldsymbol{\phi}(\mathbf{x})$. Note that the approximation (6) is equivalent with

$$\mathbf{K} \simeq \boldsymbol{\Phi}^\top \boldsymbol{\Phi}, \quad (7)$$

where \mathbf{K} is defined as the equation (3), and $\boldsymbol{\Phi}$ is defined as $\boldsymbol{\Phi} = (\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_N))$. The larger the dimension D , the higher the

approximation accuracy of the equation (6), while the larger the dimension D , the greater computational cost.

The actual function form of the feature vector $\boldsymbol{\phi}(\mathbf{x})$ depends on the kernel function. For example, in the case of the RBF kernel

$$k(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2), \quad (8)$$

which is the most famous kernel function, the vector $\boldsymbol{\phi}(\mathbf{x})$ is given by

$$\boldsymbol{\phi}(\mathbf{x}) = \cos(\mathbf{W}\mathbf{x} + \mathbf{u}), \quad (9)$$

where, the matrix $\mathbf{W} \in \mathbb{R}^{D \times M}$ is a random matrix in which each element is sampled from the normal distribution $\mathcal{N}(0, \frac{1}{4\gamma})$:

$$\mathbf{W} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{4\gamma} \mathbf{I}\right), \quad (10)$$

and the vector $\mathbf{u} \in \mathbb{R}^M$ is a random vector sampled from the uniform distribution over the range $[0, 2\pi)$:

$$\mathbf{u} \sim \mathcal{U}[0, 2\pi). \quad (11)$$

This approximation method is called random Fourier features because the feature extraction vector $\boldsymbol{\phi}(\mathbf{x})$ is derived from the Fourier transform of the kernel matrix $\mathcal{F}[k](\omega)$. Please refer to the original paper [1] for more details.

3 Gaussian Process Model and RFF

In this section, we apply RFF to Gaussian process model and theoretically confirm the effect of speeding up.

3.1 Computational complexity of Gaussian process model before applying RFF

First, let's check the computational cost required for training and inferring a normal Gaussian process model. As a premise, it is assumed that the number of training data $N \in \mathbb{Z}^+$ is sufficiently larger than the dimension $D \in \mathbb{Z}^+$ which is a hyperparameter of RFF. Here, the bottleneck of training computational cost is the calculation of the inverse matrix $(\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$ in the formulas (1) and (2). Since the size of this matrix is $N \times N$, the computational cost for the training is $O(N^3)$. Next, the bottleneck of the inference is the matrix multiplication $(\mathbf{y} - \hat{\mathbf{m}})^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$ for the expectation prediction and $\mathbf{k}(\boldsymbol{\xi}_1)^\top (\mathbf{K} - \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\boldsymbol{\xi}_2)$, for the covariance prediction whose computational costs is $O(N)$ and $O(N^2)$ respectively.

3.2 Applying RFF to expectation of prediction

Now, let's apply RFF to Gaussian process model. First of all, if you substitute the RFF approximation formula (7) into the formula of expectation of the prediction in Gaussian process (1), you'll get

$$m(\boldsymbol{\xi}) = \hat{m}(\boldsymbol{\xi}) + (\mathbf{y} - \hat{\mathbf{m}})^\top (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \boldsymbol{\phi}(\boldsymbol{\xi}), \quad (12)$$

where the matrix $\boldsymbol{\Phi} \in \mathbb{R}^{D \times N}$ is defined as $\boldsymbol{\Phi} = (\boldsymbol{\phi}(\mathbf{x}_1), \dots, \boldsymbol{\phi}(\mathbf{x}_N))$. However, this has not yet speeded up. The complexity bottleneck of the above expression (12) is still the inverse of the $N \times N$ matrix $(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I})^{-1}$.

Now we will add a bit of contrivance to the equation (12).

Lemma 1 (inversion of the matrix $\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \eta \mathbf{I}$)

Let $\boldsymbol{\Phi} \in \mathbb{R}^{D \times N}$ be a real matrix and $\sigma \in \mathbb{R}^+$ be a positive real number. Then the matrices $\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I}_N$ and $\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I}_D$ are regular matrices and the equation

$$(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I}_N)^{-1} \boldsymbol{\Phi}^\top = \boldsymbol{\Phi}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I}_D)^{-1}, \quad (13)$$

holds, where the matrices $\mathbf{I}_D \in \mathbb{R}^{D \times D}$ and $\mathbf{I}_N \in \mathbb{R}^{N \times N}$ are D and N dimensional identity matrices, respectively.

The proof of the above lemma is given at the end of this article, and let us move on to the utilization of the lemma to the equation (12). By the lemma 1, the equation (12) can be transformed as below:

$$\begin{aligned} m(\boldsymbol{\xi}) &= \hat{m}(\boldsymbol{\xi}) + (\mathbf{y} - \hat{\mathbf{m}})^\top (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \boldsymbol{\phi}(\boldsymbol{\xi}) \\ &= \hat{m}(\boldsymbol{\xi}) + (\mathbf{y} - \hat{\mathbf{m}})^\top \boldsymbol{\Phi}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})^{-1} \boldsymbol{\phi}(\boldsymbol{\xi}), \end{aligned} \quad (14)$$

provided that the two \mathbf{I} 's appearing in the above equation are of different sizes.

Clever readers would have already noticed that the bottleneck has been resolved. The inverse matrix $(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I})^{-1}$, which was the bottleneck of the expression (12), became $(\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})^{-1}$ in the expressions (14) where the size of the inverse matrix is $D \times D$. Normally, the dimension D is set sufficiently smaller than the number of training data N , therefore the inverse matrix $(\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})^{-1}$ is no longer a bottleneck of computational cost. The new bottleneck of the expression (14) is the matrix product $\boldsymbol{\Phi} \boldsymbol{\Phi}^\top$, whose computational cost is $O(ND^2)$. Therefore we've achieved a considerable speedup of the training of Gaussian process model by applying RFF because the calculational cost before RFF is $O(N^3)$.

3.3 Applying RFF to covariance of prediction

Next, we apply RFF to the covariance of the prediction (2). By substituting RFF approximation (7) into the formula of covariance prediction (2), we obtain

$$\begin{aligned} v(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2) &= k(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2) - \mathbf{k}(\boldsymbol{\xi}_1)^\top (\mathbf{K} - \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\boldsymbol{\xi}_2) \\ &= \boldsymbol{\phi}(\boldsymbol{\xi}_1)^\top \left\{ \mathbf{I} - \boldsymbol{\Phi} (\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma^2 \mathbf{I})^{-1} \boldsymbol{\Phi}^\top \right\} \boldsymbol{\phi}(\boldsymbol{\xi}_2), \end{aligned} \quad (15)$$

and by applying the lemma 1, we'll get

$$v(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2) = \boldsymbol{\phi}(\boldsymbol{\xi}_1)^\top \left\{ \mathbf{I} - \boldsymbol{\Phi} \boldsymbol{\Phi}^\top (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})^{-1} \right\} \boldsymbol{\phi}(\boldsymbol{\xi}_2). \quad (16)$$

The bottleneck of the expression (16) is, as well as the expectation prediction, the matrix inverse $(\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \sigma^2 \mathbf{I})^{-1}$ is no longer a bottleneck of the computational cost, and the new bottleneck is the matrix product $\boldsymbol{\Phi} \boldsymbol{\Phi}^\top$ whose calculation cost is $O(ND^2)$.

3.4 Training and inference process

The theoretical essence of the application of RFF to Gaussian process was mentioned in the previous subsections. In this section, let me review the procedure of training and inference of Gaussian process model after applying RFF. Each procedure is described in the algorithms 1 as pseudo-code, where \mathcal{D} is defined

as $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\sigma \in \mathbb{R}^+$ is a standard deviation of the measurement error, and ϕ is a feature extraction function used in the RFF approximation.

In the algorithm 1, the function TRAINING_GP_WITH_RFF is designed to take the training data \mathcal{D} and the standard deviation σ as arguments, and return a vector \mathbf{a} and a matrix \mathbf{C} that will be used in the inference of expectation and covariance, respectively. Also, the function INFERENCE_GP_WITH_RFF is designed to take the inference target data ξ and the training result \mathbf{a} and \mathbf{C} , and return predicted expectation μ and predicted variance v . In the algorithm 1, we computed only the variance of prediction, however, if you want to compute the covariance of two input data ξ_1 and ξ_2 , then please compute $\phi(\xi_1)^\top \mathbf{C} \phi(\xi_2)$. Also, note that the prior distribution of Gaussian process model is set to 0 for the sake of simplicity in the algorithm 1.

Algorithm 1 Training of the GP model after RFF

```

1: function TRAINING_GP_WITH_RFF( $\mathcal{D}, \sigma$ )
2:    $\mathbf{y} \leftarrow (y_1, y_2, \dots, y_N)^\top$ 
3:    $\Phi \leftarrow (\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N))$ 
4:    $\mathbf{P} \leftarrow \Phi \Phi^\top$ 
5:    $\mathbf{Q} \leftarrow (\mathbf{P} + \sigma^2 \mathbf{I})^{-1}$ 
6:    $\mathbf{a} \leftarrow \mathbf{y}^\top \mathbf{Q}$ 
7:    $\mathbf{C} \leftarrow \mathbf{I} - \mathbf{P} \mathbf{Q}$ 
8:   return ( $\mathbf{a}, \mathbf{C}$ )
9: end function
10:
11: function INFERENCE_GP_WITH_RFF( $\xi, \mathbf{a}, \mathbf{C}$ )
12:    $\mathbf{z} \leftarrow \phi(\xi)$ 
13:    $\mu \leftarrow \mathbf{a}^\top \mathbf{z}$ 
14:    $v \leftarrow \mathbf{z}^\top \mathbf{C} \mathbf{z}$ 
15:   return ( $\mu, v$ )
16: end function

```

3.5 Computational cost of training and inference

Finally, let me summarize the computational cost of the training and inference of Gaussian process model with random Fourier features. In this section, we assume that the number of training data N is sufficiently larger than the RFF dimension D . Also the computational cost of the feature extraction function ϕ is ignored because it depends on the kernel function itself.

As you can see the algorithm 1, the bottleneck of the training process is the matrix multiplication $\Phi \Phi^\top$ whose computational cost is $O(ND^2)$. Note that the matrix inversion $(\mathbf{P} + \sigma^2 \mathbf{I})^{-1}$ where $\mathbf{P} = \Phi \Phi^\top$ is not a bottleneck of the training process under the assumption $N \gg D$, because its computational cost is $O(D^3)$.

As for the inference cost, the bottleneck is $\mathbf{a}^\top \mathbf{z}$ for the expectation prediction, and $\mathbf{z}^\top \mathbf{C} \mathbf{z}$ for the covariance prediction that have the computational cost $O(D)$ and $O(D^2)$ respectively. The computational cost before and after applying RFF is summarized in the table 1.

Table 1: Computational cost of the GP model before/after RFF

	Training	Inference
Before RFF	$O(N^3)$	$O(N^2)$
After RFF	$O(ND^2)$	$O(D^2)$

A Appendix 1: Another Approach

This section provides another approach to reduce the computational cost of the equation (12). This approach has almost the same computational cost as the previous approach, however, the previous approach is a bit more simple and beautiful^{†1} than this approach. Therefore the readers don't have to pay much attention to this section.

A.1 Applying RFF to expectation of prediction

First, let us introduce the *matrix inversion lemma* (it's also referred to as the *binominal inverse lemma*) which is a useful formula for the expansion of a matrix inverse.

Lemma 2 (Matrix Inversion Lemma)

Let $\mathbf{A} \in \mathbb{R}^{N \times N}$, $\mathbf{B} \in \mathbb{R}^{N \times M}$, $\mathbf{C} \in \mathbb{R}^{M \times N}$, and $\mathbf{D} \in \mathbb{R}^{M \times M}$ be real matrices. Then the equation

$$(\mathbf{A} + \mathbf{B} \mathbf{D} \mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D}^{-1} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C} \mathbf{A}^{-1} \quad (17)$$

holds, where the matrices \mathbf{A} and \mathbf{D} are regular matrices.

The proof of the matrix inversion lemma is given at the end of this article, and let us move on to the utilization of the lemma to the equation (12).

By replacing $\mathbf{A} = \sigma^2 \mathbf{I}$, $\mathbf{B} = \Phi^\top$, $\mathbf{C} = \Phi$, and $\mathbf{D} = \mathbf{I}$ on the equation (17), we obtain the following equation:

$$(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} = \frac{1}{\sigma^2} \left(\mathbf{I} - \Phi^\top (\Phi \Phi^\top + \sigma^2 \mathbf{I})^{-1} \Phi \right), \quad (18)$$

where $\mathbf{P} = \Phi \Phi^\top \in \mathbb{R}^{D \times D}$. Then multiply Φ from the right to the above equation (18), and we get

$$(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} \Phi^\top = \frac{1}{\sigma^2} \Phi^\top \left(\mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P} \right). \quad (19)$$

Therefore, the expression (12) can be written as

$$m(\xi) = \hat{m}(\xi) + \frac{1}{\sigma^2} (\mathbf{y} - \hat{\mathbf{m}})^\top \Phi^\top \mathbf{S} \phi(\xi), \quad (20)$$

where

$$\mathbf{S} = \mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P}. \quad (21)$$

The same as the previous approach, the bottleneck has been resolved now. The inverse matrix $(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1}$, which was the bottleneck of the expression (12), became $(\mathbf{P} + \sigma^2 \mathbf{I})^{-1}$ in the expressions (20) and (21) where the size of the inverse matrix is $D \times D$. Normally, the RFF dimension D is set sufficiently smaller than the number of training data N , therefore the inverse matrix $(\mathbf{P} + \sigma^2 \mathbf{I})^{-1}$ is no longer a bottleneck of computational cost. The bottleneck of the expressions (20) and (21) is the matrix product $\mathbf{P} = \Phi \Phi^\top$, whose computational cost is $O(ND^2)$. Therefore we've achieved a considerable speedup of the training of Gaussian process model by applying RFF because the calculational cost before RFF is $O(N^3)$.

A.2 Applying RFF to covariance of prediction

Next, we apply RFF to the covariance of the prediction (2). By substituting RFF approximation (19) to the expression (2), we

^{†1}The previous approach has high similarity with the linear regression, however, the approach mentioned in this section doesn't have.

obtain

$$\begin{aligned} v(\xi_1, \xi_2) &= \phi(\xi_1)^\top \phi(\xi_2) - \frac{1}{\sigma^2} \phi(\xi_1)^\top \mathbf{P} \mathbf{S} \phi(\xi_2) \\ &= \phi(\xi_1)^\top \left(\mathbf{I} - \frac{1}{\sigma^2} \mathbf{P} \mathbf{S} \right) \phi(\xi_2), \end{aligned} \quad (22)$$

The bottleneck of the expression (16) is, as well as the expectation of the prediction, the matrix product $\mathbf{P} = \Phi \Phi^\top$ whose calculation cost is $O(ND^2)$.

The procedure of training and inference of Gaussian process model after applying RFF is described in the algorithms 2 as pseudo-code. Note that the prior distribution of Gaussian process model is set to 0 for the sake of simplicity in the algorithms 2.

Algorithm 2 Training of the GP model after RFF (2)

```

1: function TRAINING_GP_WITH_RFF( $\mathcal{D}, \sigma$ )
2:    $\mathbf{y} \leftarrow (y_1, \dots, y_N)^\top$ 
3:    $\Phi \leftarrow (\phi(x_1), \dots, \phi(x_N))$ 
4:    $\mathbf{P} \leftarrow \Phi \Phi^\top$ 
5:    $\mathbf{S} \leftarrow \mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P}$ 
6:    $\mathbf{a} \leftarrow \frac{1}{\sigma^2} \mathbf{y}^\top \Phi^\top \mathbf{S}$ 
7:    $\mathbf{C} \leftarrow \mathbf{I} - \frac{1}{\sigma^2} \mathbf{P} \mathbf{S}$ 
8:   return ( $\mathbf{a}, \mathbf{C}$ )
9: end function
10:
11: function INFERENCE_GP_WITH_RFF( $\mathcal{D}, \sigma$ )
12:    $\mathbf{z} \leftarrow \phi(\xi)$ 
13:    $\mu \leftarrow \mathbf{a}^\top \mathbf{z}$ 
14:    $v \leftarrow \mathbf{z}^\top \mathbf{C} \mathbf{z}$ 
15:   return ( $\mu, v$ )
16: end function

```

B Appendix 2: Proofs

B.1 Proof of the lemma 1

The lemma 1 is reprinted and proved.

Lemma 1 (inversion of the matrix $\Phi^\top \Phi + \eta \mathbf{I}$)

Let $\Phi \in \mathbb{R}^{D \times N}$ be a real matrix and η be a positive real number. Then the matrices $\Phi^\top \Phi + \eta \mathbf{I}_N$ and $\Phi \Phi^\top + \eta \mathbf{I}_D$ are regular matrices and the equation

$$\left(\Phi^\top \Phi + \eta \mathbf{I}_N \right)^{-1} \Phi^\top = \Phi^\top \left(\Phi \Phi^\top + \eta \mathbf{I}_D \right)^{-1} \quad (23)$$

holds, where the matrices \mathbf{I}_D and \mathbf{I}_N are D dimensional and N dimensional identity matrices, respectively.

Proof: First, the matrices $\Phi^\top \Phi + \eta \mathbf{I}_N$ and $\Phi \Phi^\top + \eta \mathbf{I}_D$ are positive definite, because

$$\begin{aligned} \mathbf{x}^\top \left(\Phi^\top \Phi + \eta \mathbf{I}_N \right) \mathbf{x} &= \|\Phi \mathbf{x}\|^2 + \eta \|\mathbf{x}\|^2 > 0, \\ \hat{\mathbf{x}}^\top \left(\Phi \Phi^\top + \eta \mathbf{I}_D \right) \hat{\mathbf{x}} &= \|\Phi^\top \hat{\mathbf{x}}\|^2 + \eta \|\hat{\mathbf{x}}\|^2 > 0, \end{aligned}$$

holds for any non-zero real vectors $\mathbf{x} \in \mathbb{R}^N$ and $\hat{\mathbf{x}} \in \mathbb{R}^D$. In general, the determinant of a matrix is equal to the product of its eigenvalues. Also, the eigenvalues of a positive definite matrix are always greater than zero. Therefore, the determinant of a positive definite matrix is greater than zero. Hence the matrices $\Phi^\top \Phi + \eta \mathbf{I}_N$

and $\Phi \Phi^\top + \eta \mathbf{I}_D$ are regular matrices, that is, these matrices are invertible.

Next, it is obvious that the equation

$$\Phi^\top \left(\Phi \Phi^\top + \eta \mathbf{I}_D \right) = \left(\Phi^\top \Phi + \eta \mathbf{I}_N \right) \Phi^\top,$$

holds, because both are equivalent to $\Phi^\top \Phi \Phi^\top + \eta \Phi^\top$. You'll get the following equation by multiplying the inverse matrix $\left(\Phi^\top \Phi + \eta \mathbf{I}_N \right)^{-1}$ to the above equation from the left:

$$\left(\Phi^\top \Phi + \eta \mathbf{I}_N \right)^{-1} \Phi^\top \left(\Phi \Phi^\top + \eta \mathbf{I}_D \right) = \Phi^\top.$$

Similarly, by multiplying the inverse matrix $\left(\Phi \Phi^\top + \eta \mathbf{I}_D \right)^{-1}$ to the above equation from the right, you'll get

$$\left(\Phi^\top \Phi + \eta \mathbf{I}_N \right)^{-1} \Phi^\top = \Phi^\top \left(\Phi \Phi^\top + \eta \mathbf{I}_D \right)^{-1}.$$

B.2 Proof of the matrix inversion lemma

The matrix inversion lemma is reprinted and proved.

Lemma 2 (Matrix Inversion Lemma)

Let $\mathbf{A} \in \mathbb{R}^{N \times N}$, $\mathbf{B} \in \mathbb{R}^{N \times M}$, $\mathbf{C} \in \mathbb{R}^{M \times N}$, and $\mathbf{D} \in \mathbb{R}^{M \times M}$ be real matrices. Then the equation

$$\left(\mathbf{A} + \mathbf{B} \mathbf{D} \mathbf{C} \right)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} \left(\mathbf{D}^{-1} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \right)^{-1} \mathbf{C} \mathbf{A}^{-1} \quad (24)$$

holds, where the matrices \mathbf{A} and \mathbf{D} are regular matrices.

Proof: The following equation holds:

$$\begin{aligned} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} &= \begin{pmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1} \mathbf{B} \mathbf{S} \mathbf{C} \mathbf{A}^{-1} & -\mathbf{A}^{-1} \mathbf{B} \mathbf{S} \\ -\mathbf{S} \mathbf{C} \mathbf{A}^{-1} & \mathbf{S} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{T} & -\mathbf{T} \mathbf{B} \mathbf{D}^{-1} \\ -\mathbf{D}^{-1} \mathbf{C} \mathbf{T} & \mathbf{D}^{-1} + \mathbf{D}^{-1} \mathbf{C} \mathbf{T} \mathbf{B} \mathbf{D}^{-1} \end{pmatrix}, \end{aligned}$$

where

$$\mathbf{T} = \left(\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \right)^{-1}, \quad (25)$$

$$\mathbf{S} = \left(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} \right)^{-1}. \quad (26)$$

It is easy to verify the above equation from a direct calculation. By comparing the corresponding parts of the above block matrix, we get

$$\mathbf{T} = \mathbf{A}^{-1} + \mathbf{A}^{-1} \mathbf{B} \mathbf{S} \mathbf{C} \mathbf{A}^{-1}, \quad (27)$$

$$\mathbf{S} = \mathbf{D}^{-1} + \mathbf{D}^{-1} \mathbf{C} \mathbf{T} \mathbf{B} \mathbf{D}^{-1}, \quad (28)$$

$$-\mathbf{A}^{-1} \mathbf{B} \mathbf{S} = -\mathbf{T} \mathbf{B} \mathbf{D}^{-1}, \quad (29)$$

$$-\mathbf{S} \mathbf{C} \mathbf{A}^{-1} = -\mathbf{D}^{-1} \mathbf{C} \mathbf{T}, \quad (30)$$

By replacing with

$$\mathbf{A} \rightarrow \mathbf{D}^{-1}, \quad \mathbf{B} \rightarrow -\mathbf{C}, \quad \mathbf{C} \rightarrow \mathbf{B}, \quad \mathbf{D} \rightarrow \mathbf{A},$$

in the equation (27), we get the formula to be proved. ■

References

- [1] A. Rahimi and B. Recht, “Random Features for Large-Scale Kernel Machines”, Neural Information Processing Systems, 2007.
- [2] C. Rasmussen and C. Williams, “Gaussian Processes for Machine Learning”, MIT Press, 2006.
- [3] <https://github.com/tiskw/mathematical-articles>