

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 the 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, like 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 the 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 Gaussian process model is $O(N^3)$, and 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 the kernel method, in other words, the kernel support vector machine also has the same problem.

One of the methods to speed up the kernel method is random Fourier features [1] (hereinafter abbreviated as RFF). This method can significantly reduce the computational cost while keeping the flexibility of the kernel method by approximating the kernel function as the 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 the kernel method, 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 articles 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 this repository¹.

1 Gaussian Process Model Revisited

This section gives an overview of Gaussian process model. Unfortunately, this section can not cover details such as the formulation and derivation of Gaussian process models due to 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 a 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 normal distribution. If the test data is $\xi \in \mathbb{R}^M$, the expectation of the prediction is given by:

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

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, $\hat{m}(\xi)$ is the prior distribution of the prediction, and $\hat{\mathbf{m}} = (\hat{m}(\mathbf{x}_1), \dots, \hat{m}(\mathbf{x}_N))^\top$ 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 $\hat{m}(\cdot) = 0$ and $\hat{\mathbf{m}} = \mathbf{0}$.

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

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 the kernel function. In RFF, the kernel function can be approximated as

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

where $\phi(\mathbf{x}) \in \mathbb{R}^D$ is a feature vector extracted from the data \mathbf{x} and $D \in \mathbb{Z}^+$ is the dimension of $\phi(\mathbf{x})$. The larger the dimension D , the higher the approximation accuracy of the equation (6), while the larger the dimension D , the greater computational cost.

¹<https://github.com/tiskw/mathematical-articles>

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

$$k(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|^2) \quad (7)$$

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

$$\phi(\mathbf{x}) = \cos(\mathbf{W}\mathbf{x} + \mathbf{u}) \quad (8)$$

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 $\mathcal{N}(0, \frac{1}{4\gamma})$:

$$\mathbf{W} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{4\gamma} \mathbf{I}\right), \quad (9)$$

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

$$\mathbf{u} \sim \mathcal{U}[0, 2\pi). \quad (10)$$

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 $M \in \mathbb{Z}^+$ of the input vector and dimension $D \in \mathbb{Z}^+$ which is a hyperparameter of RFE. Here, the bottleneck of training computational cost is obviously 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 training is $O(N^3)$.

Next, the bottleneck of the inference is matrix multiplications $(\mathbf{y} - \hat{\mathbf{m}})^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$ or $\mathbf{k}(\xi_1)^\top (\mathbf{K} - \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\xi_2)$, and either of these computational cost is $O(N^2)$.

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 (2) into the formula of expectation of the prediction in Gaussian process (1), you'll get

$$m(\xi) = \hat{m}(\xi) + (\mathbf{y} - \hat{\mathbf{m}})^\top (\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} \Phi^\top \phi(\xi), \quad (11)$$

where the matrix $\Phi \in \mathbb{R}^{D \times N}$ is defined as $\Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N))$. However, this has not yet speeded up. The complexity bottleneck of the above expression (11) is still the inverse of the $N \times N$ matrix $(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1}$.

Now we will add a bit of contrivance to the equation (11). At first, let us introduce the *matrix inversion lemma* (it's also referred as *binominal inverse lemma*) which is a useful formula for expansion of matrix inverse.

Theorem 3.1 (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

$$(\mathbf{A} + \mathbf{B}\mathbf{D}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D}^{-1} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} \quad (12)$$

holds, where the matrix \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 (11).

By replacing $\mathbf{A} = \sigma^2 \mathbf{I}$, $\mathbf{B} = \Phi^\top$, $\mathbf{C} = \Phi$, and $\mathbf{D} = \mathbf{I}$ on the equation (12), we obtain the following equation:

$$(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} = \frac{1}{\sigma^2} \left(\mathbf{I} - \Phi^\top (\Phi \Phi^\top + \sigma^2 \mathbf{I})^{-1} \Phi \right), \quad (13)$$

where $\mathbf{P} = \Phi \Phi^\top \in \mathbb{R}^{D \times D}$. Then multiply Φ from the right to the above equation (13), we get

$$(\Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} \Phi^\top = \frac{1}{\sigma^2} \Phi^\top (\mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P}). \quad (14)$$

Therefore, the expression (11) can be written as

$$m(\xi) = \hat{m}(\xi) + \frac{1}{\sigma^2} (\mathbf{y} - \hat{\mathbf{m}})^\top \Phi^\top \mathbf{S} \phi(\xi), \quad (15)$$

where

$$\mathbf{S} = \mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P}. \quad (16)$$

Clever readers would have already noticed that the bottleneck has been resolved. The inverse matrix $(\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$, which was the bottleneck of the expression (11), became $(\mathbf{P} + \sigma^2 \mathbf{I})^{-1}$ in the expressions (15) and (16) 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 (15) and (16) 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 $O(N^3)$.

3.3 Applying RFF to covariance of prediction

Next, we apply RFF to the covariance of the prediction (2). By substitute RFF approximation (14) to the expression (2), we obtain

$$\begin{aligned} v(\xi_1, \xi_2) &= \phi(\xi_1)^\top \phi(\xi_2) - \frac{1}{\sigma^2} \phi(\xi_1)^\top \mathbf{P} \mathbf{S} \phi(\xi_2) \\ &= \phi(\xi_1)^\top \left(\mathbf{I} - \frac{1}{\sigma^2} \mathbf{P} \mathbf{S} \right) \phi(\xi_2), \end{aligned} \quad (17)$$

The bottleneck of the expression (17) is, as the same 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 algorithm 1 and 2 as pseudo code. Note that the prior distribution of Gaussian process model is set to 0 for the sake of simplicity in Algorithm 1 and 2.

Finally, the calculational cost after applying RFF is summarized in the table 1, where $N \in \mathbb{Z}^+$ is the number of training data and $D \in \mathbb{Z}^+$ is the dimension of RFF.

Algorithm 1: Training of the GP model after RFF

Data: $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\sigma \in \mathbb{R}^+$
Result: $\mathbf{c}_m \in \mathbb{R}^D$, $\mathbf{C}_v \in \mathbb{R}^{D \times D}$
 $\mathbf{y} \leftarrow (y_1, \dots, y_N)^\top$
 $\Phi \leftarrow (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N))$
 $\mathbf{P} \leftarrow \Phi \Phi^\top$
 $\mathbf{S} \leftarrow \mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1} \mathbf{P}$
 $\mathbf{c}_m \leftarrow \frac{1}{\sigma^2} \mathbf{y}^\top \Phi^\top \mathbf{S}$ /* Cache for expectation */
 $\mathbf{C}_v \leftarrow \mathbf{I} - \frac{1}{\sigma^2} \mathbf{P} \mathbf{S}$ /* Cache for covariance */

Algorithm 2: Inference of the GP model after RFF

Data: $\xi \in \mathbb{R}^M$, $\mathbf{c}_m \in \mathbb{R}^D$, $\mathbf{C}_v \in \mathbb{R}^{D \times D}$
Result: $\mu \in \mathbb{R}$, $\eta \in \mathbb{R}$
 $\mathbf{z} \leftarrow \phi(\xi)$
 $\mu \leftarrow \mathbf{c}_m^\top \mathbf{z}$ /* Inference of expectation */
 $\eta \leftarrow \mathbf{z}^\top \mathbf{C}_v \mathbf{z}$ /* Inference of covariance */

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

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

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.

A Appendix: Proofs**A.1 Proof of matrix inversion lemma**

The matrix inversion lemma is reprinted and proved.

Theorem A.1 (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

$$(\mathbf{A} + \mathbf{B} \mathbf{D} \mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{D}^{-1} + \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{C} \mathbf{A}^{-1} \quad (18)$$

holds, where the matrix \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

$$\mathbf{T} = (\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B})^{-1}, \quad (19)$$

$$\mathbf{S} = (\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1}. \quad (20)$$

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

$$\mathbf{T} = \mathbf{A}^{-1} + \mathbf{A}^{-1} \mathbf{B} \mathbf{S} \mathbf{C} \mathbf{A}^{-1}, \quad (21)$$

$$\mathbf{S} = \mathbf{D}^{-1} + \mathbf{D}^{-1} \mathbf{C} \mathbf{T} \mathbf{B} \mathbf{D}^{-1}, \quad (22)$$

$$-\mathbf{A}^{-1} \mathbf{B} \mathbf{S} = -\mathbf{T} \mathbf{B} \mathbf{D}^{-1}, \quad (23)$$

$$-\mathbf{S} \mathbf{C} \mathbf{A}^{-1} = -\mathbf{D}^{-1} \mathbf{C} \mathbf{T}, \quad (24)$$

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 (21), we get the formula to be proved. \blacksquare