

Gaussian Process Bootstrapping Layer

Tetsuya Ishikawa

October 10, 2021

1 Theoretical details

The Gaussian process with random Fourier features can be regarded as a variant of a fully connected layer. We can design a new fully connected layer that can compute variance of its intermediate features.

1.1 Gaussian process with random Fourier features

Let $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ be a training dataset where \mathbf{x}_n is a n -th data and y_n is a label of n -th data. The Gaussian process with random Fourier features can be formulated as follows:

$$\boldsymbol{\phi}_{\mathbf{x}_n} = f(\mathbf{W}\mathbf{x}_n), \quad (1)$$

$$\mathbf{P} = \sum_{n=1}^N \boldsymbol{\phi}_{\mathbf{x}_n} \boldsymbol{\phi}_{\mathbf{x}_n}^\top, \quad (2)$$

$$m(\mathbf{x}_i) = \frac{1}{\sigma^2} \mathbf{y}^\top \boldsymbol{\Phi}^\top (\mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1}) \boldsymbol{\phi}_{\mathbf{x}_i}, \quad (3)$$

$$v(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}_{\mathbf{x}_i}^\top \left\{ \mathbf{I} - \frac{1}{\sigma^2} \mathbf{P} (\mathbf{I} - (\mathbf{P} + \sigma^2 \mathbf{I})^{-1}) \right\} \boldsymbol{\phi}_{\mathbf{x}_j}, \quad (4)$$

where f is a non-linear function and \mathbf{W} is a random matrix derived from a kernel function of the Gaussian process.

1.2 Analogy to a fully connected layer

The formula $f(\mathbf{W}\mathbf{x})$ looks like a fully connected layer of neural network with activation function. For example, if the kernel function is RBF kernel $k(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\|\mathbf{x}_1 - \mathbf{x}_2\|^2)$, then $\boldsymbol{\phi}_{\mathbf{x}_n}$ can be written as

$$\boldsymbol{\phi}_{\mathbf{x}_n} = \begin{pmatrix} \cos \mathbf{W}\mathbf{x}_n \\ \sin \mathbf{W}\mathbf{x}_n \end{pmatrix}, \quad \mathbf{W} \sim \mathcal{N}(0, \sigma^2). \quad (5)$$

In this case, $\mathbf{W}\mathbf{x}_n$ corresponds to a fully connected layer and the circular functions correspond to an activation function.

On the other hand, $v(\mathbf{x}_n, \mathbf{x}_n)$ correspond to the variance of input data \mathbf{x}_n . Notable point is that the $v(\mathbf{x}_n, \mathbf{x}_n)$ is not depend on the label y_n , in other words, variance $v(\mathbf{x}_n, \mathbf{x}_n)$ is the same for any label. See the equation (4).

1.3 Gaussian process bootstrapping layer

Because of the independence of the variance and $v(\mathbf{x}_n, \mathbf{x}_n)$ and the label y_n mentioned in the previous subsection, we can replace the expectation prediction part as a identity function (theoretically, this situation correspond that $y_n = \boldsymbol{\phi}_{\mathbf{x}_n}$). See the figure 1. Therefore, we've got a new layer that can predict variance of intermediate features by replacing the random Fourier features as a fully connected layer, and expectation prediction as a identity function.

Gaussian process bootstrapping layer is a layer to add noises to the intermediate features where the variance of the noises is the variance of the intermediate features.

1.4 Psuedo code of GPB layer

The algorithm 1 is the pseudo code of the GPB layer.

References

- [1] C. Rasmussen and C. Williams, "Gaussian Processes for Machine Learning", MIT Press, 2006.
- [2] A. Rahimi and B. Recht, "Random Features for Large-Scale Kernel Machines", NIPS, 2007.

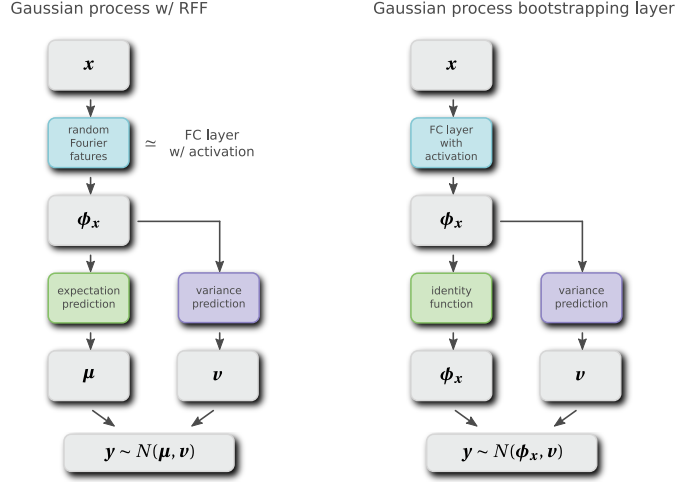


Figure 1: Illustration of the analogy between GP w/ RFF and GPB layer

Algorithm 1 Gaussian process bootstrapping layer

Input

X : tensor with shape (N, C)

Output

Y : tensor with shape (N, C)

Hyperparameters

σ : standard deviation of measurement error

α : coefficient of exponential moving average

s : number of steps to skip bootstrapping

function GAUSSIAN PROCESS BOOTSTRAPPING LAYER(X, α, σ)

Update matrix P with exponential moving average.

$$P = \alpha X^T X + (1 - \alpha)P$$

Compute matrix M .

$$M = I - \frac{1}{\sigma^2} P (I - (P + \sigma^2 * I)^{-1} P)$$

Compute variance $v[n]$.

for n **in** $[0, N)$:

$$v[n] = X[n, :]^T M X[n, :]$$

Add perturbation to the input tensor X .

for n **in** $[0, N)$:

$$Y[n, :] = X[n, :] + \sqrt{v[n]} \text{ (sampling from normal distribution with shape } (1, C) \text{)}$$

end function
