

# **Exercise of Supervised Learning: Gaussian Processes Part 2**

Yawei Li

`yawei.li@stat.uni-muenchen.de`

February 4, 2025

# Exercise 1: Gaussian Posterior Process

Assume your data follows the following law:

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

with  $\mathbf{f} = f(\mathbf{x}) \in \mathbb{R}^n$  being a realization of a Gaussian process (GP), for which we a priori assume

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')).$$

$\mathbf{x}$  here only consists of 1 feature that is observed for  $n$  data points.

(a) Derive / define the prior distribution of  $\mathbf{f}$ .

# 1 (a): Prior Distribution of $\mathbf{f}$

- ▶  $\mathbf{f} \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$ .
- ▶  $\mathbf{m} = m(\mathbf{x})$ .
- ▶  $\mathbf{K}_{ij} = k(x^{(i)}, x^{(j)})$ .
- ▶ NB: Note the (in-)finite Gaussian property of a GP.

## Exercise 1 (b)

(b) Derive the posterior distribution of  $\mathbf{f}|\mathbf{y}$ .

# 1 (b): Likelihood and Prior

The likelihood is:

$$p(\mathbf{y}|\mathbf{f}) \propto \exp \left( -\frac{1}{2}(\mathbf{y} - \mathbf{f})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{f}) \right).$$

The prior is:

$$p(\mathbf{f}) \propto \exp \left( -\frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m}) \right).$$

# 1 (b): Posterior

So the posterior is:

$$\begin{aligned} p(\mathbf{f}|\mathbf{y}) &\propto p(\mathbf{y}|\mathbf{f}) \cdot p(\mathbf{f}) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{f})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{f})\right) \cdot \exp\left(-\frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m})\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{f}^\top (\sigma^2 \mathbf{I})^{-1} \mathbf{f} - 2\mathbf{f}^\top (\sigma^2 \mathbf{I})^{-1} \mathbf{y} + \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f} - 2\mathbf{f}^\top \mathbf{K}^{-1} \mathbf{m})\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{f}^\top ((\sigma^2 \mathbf{I})^{-1} + \mathbf{K}^{-1}) \mathbf{f} - 2 \cdot \mathbf{f}^\top ((\sigma^2 \mathbf{I})^{-1} \mathbf{y} + \mathbf{K}^{-1} \mathbf{m}))\right) \\ &\propto \exp\left(-\frac{1}{2} \mathbf{f}^\top \underbrace{((\sigma^2 \mathbf{I})^{-1} + \mathbf{K}^{-1})}_{=:\mathbf{K}_{\text{post}}^{-1}} \mathbf{f} + \mathbf{f}^\top \underbrace{((\sigma^2 \mathbf{I})^{-1} \mathbf{y} + \mathbf{K}^{-1} \mathbf{m})}_{=:\tilde{\mathbf{f}}}\right) \\ &\propto \exp\left(-\frac{1}{2} \mathbf{f}^\top \mathbf{K}_{\text{post}}^{-1} \mathbf{f} + \mathbf{f}^\top \tilde{\mathbf{f}}\right) \end{aligned}$$

# 1 (b): Complete the Square

$$p(\mathbf{f}|\mathbf{y}) \propto \exp\left(-\frac{1}{2}\mathbf{f}^\top \mathbf{K}_{\text{post}}^{-1} \mathbf{f} + \mathbf{f}^\top \tilde{\mathbf{f}}\right)$$

Recall the technique of **completing the square**:

- Scalar:  $ax^2 + bx + c = a(x + \frac{b}{2a})^2 + (c - \frac{b^2}{4a})$
- Matrix / vector:  $\mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{x}^\top \mathbf{b} + c = (\mathbf{x} + \frac{1}{2}\mathbf{A}^{-1}\mathbf{b})^\top \mathbf{A}(\mathbf{x} + \frac{1}{2}\mathbf{A}^{-1}\mathbf{b}) + (c - \frac{1}{4}\mathbf{b}^\top \mathbf{A}^{-1}\mathbf{b})$  for a **symmetric** matrix  $\mathbf{A}$ .

In our case:  $\mathbf{A} = -\frac{1}{2}\mathbf{K}_{\text{post}}^{-1}$ , and  $\mathbf{b} = \tilde{\mathbf{f}}$ , and  $c = 0$  or an arbitrary const. So,

$$\mathbf{x} + \frac{1}{2}\mathbf{A}^{-1}\mathbf{b} = \mathbf{f} + \frac{1}{2}\left(-\frac{1}{2}\mathbf{K}_{\text{post}}^{-1}\right)^{-1}\tilde{\mathbf{f}} = \mathbf{f} - \mathbf{K}_{\text{post}}\tilde{\mathbf{f}}$$

and we omit  $c - \frac{1}{4}\mathbf{b}^\top \mathbf{A}^{-1}\mathbf{A}\mathbf{b}$  because

- It is a constant w.r.t.  $\mathbf{f}$ .
- the exp and  $\propto$  operators.

## 1 (b): Complete the Square

$$\begin{aligned} p(\mathbf{f}|\mathbf{y}) &\propto \exp \left( (\mathbf{f} - \mathbf{K}_{\text{post}} \tilde{\mathbf{f}})^\top \left( -\frac{1}{2} \mathbf{K}_{\text{post}}^{-1} \right)^{-1} (\mathbf{f} - \underbrace{\mathbf{K}_{\text{post}} \tilde{\mathbf{f}}}_{:=\mathbf{f}_{\text{post}}}) \right) \\ &\propto \exp \left( (\mathbf{f} - \mathbf{f}_{\text{post}})^\top \left( -\frac{1}{2} \mathbf{K}_{\text{post}}^{-1} \right)^{-1} (\mathbf{f} - \mathbf{f}_{\text{post}}) \right) \end{aligned}$$

Hence,

$$\mathbf{f}|\mathbf{y} \sim \mathcal{N}(\mathbf{f}_{\text{post}}, \mathbf{K}_{\text{post}}).$$



## Exercise 1 (c)

(c) Derive the posterior predictive distribution  $y_* | x_*, \mathbf{x}, \mathbf{y}$  for a new sample  $x_*$  from the sample data-generating process.

# 1 (c): Derive Predictive Posterior from Joint Distribution

Naïvely, we can compute

$$p(y_*|x_*, \mathbf{x}, \mathbf{y}) = \int p(y_*|x_*, \mathbf{x}, \mathbf{y}, \mathbf{f}) \cdot p(\mathbf{f}|\mathbf{y}, \mathbf{x}) d\mathbf{f}.$$

This is cumbersome. Alternative, we can use the fact that

$$\text{if } \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ab}^\top & \Sigma_{bb} \end{pmatrix} \right), \text{ then } \mathbf{a}|\mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_{a|b}, \Sigma_{a|b})$$

where

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{b} - \boldsymbol{\mu}_b) \\ \Sigma_{a|b} &= \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \end{aligned}$$

# 1 (c) Derive Predictive Posterior from Joint Distribution

The joint distribution of  $(\mathbf{y}, y_*)$ :

**Note:** here is  $\mathbf{y}$  instead of  $\mathbf{f}$ , and  $\mathbf{y} = \mathbf{f} + \epsilon$ . So we have  $\sigma^2$  in the in the cov. matrix.

$$\begin{pmatrix} \mathbf{y} \\ y_* \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbf{m} \\ m_* \end{pmatrix}, \begin{pmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{K}_* \\ \mathbf{K}_*^\top & \mathbf{K}_{**} + \sigma^2 \end{pmatrix} \right),$$

Therefore, the conditional distribution  $y_* | x_*, \mathbf{x}, \mathbf{y}$  is also a Gaussian:

$$y_* | x_*, \mathbf{x}, \mathbf{y} \sim \mathcal{N}(m_* + \mathbf{K}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}), \mathbf{K}_{**} + \sigma^2 - \mathbf{K}_*^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_*).$$

## Exercise 1 (d)

Implement the GP with squared exponential kernel, zero mean function and  $\ell = 1$  from scratch for  $n = 2$  observations  $(\mathbf{y}, \mathbf{x})$ . Do this as efficiently as possible by explicitly calculating all expensive computations by hand. Do the same for the posterior predictive distribution of  $y_*$ . Test your implementation using simulated data.

**Show the standard solution.**