

# Gaussian Processes

## Advanced Statistical Inference

Simone Rossi

### Multivariate Gaussian: Marginalization and Conditioning

Let  $\mathbf{z} = (z_1, z_2, z_3)^\top \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with

$$\boldsymbol{\mu} = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} 4 & 2 & 0 \\ 2 & 3 & 1 \\ 0 & 1 & 2 \end{pmatrix}.$$

Partition the vector as  $\mathbf{z} = (\mathbf{z}_a^\top, z_b)^\top$  where  $\mathbf{z}_a = (z_1, z_2)^\top$  and  $z_b = z_3$ , so that

$$\boldsymbol{\mu}_a = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad \mu_b = 0, \quad \boldsymbol{\Sigma}_{aa} = \begin{pmatrix} 4 & 2 \\ 2 & 3 \end{pmatrix}, \quad \boldsymbol{\Sigma}_{ab} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \Sigma_{bb} = 2.$$

1. **Marginalization.** State the general rule for the marginal  $p(\mathbf{z}_a)$  of a partitioned Gaussian. Apply it to compute  $p(\mathbf{z}_a) = p(z_1, z_2)$  explicitly, giving its mean vector and covariance matrix. Analyze all the dimensions of the matrices and vectors involved to verify that they are consistent.
2. **Conditioning.** State the general formula for the conditional  $p(\mathbf{z}_a | z_b)$ :

$$p(\mathbf{z}_a | z_b) = \mathcal{N} \left( \underbrace{\boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \Sigma_{bb}^{-1} (z_b - \mu_b)}_{\boldsymbol{\mu}_{a|b}}, \underbrace{\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \Sigma_{bb}^{-1} \boldsymbol{\Sigma}_{ab}^\top}_{\boldsymbol{\Sigma}_{a|b}} \right).$$

For the observation  $z_b = 2$ , compute numerically:

- the conditional mean  $\boldsymbol{\mu}_{a|b}$ ;
  - the conditional covariance  $\boldsymbol{\Sigma}_{a|b}$ .
3. **Interpretation.** Compare  $\boldsymbol{\Sigma}_{aa}$  and  $\boldsymbol{\Sigma}_{a|b}$ . Which is “larger”? What does this tell you about the effect of observing  $z_b$  on our uncertainty about  $\mathbf{z}_a$ ? Which component ( $z_1$  or  $z_2$ ) is affected more, and why?

## From Parametric to Non-Parametric Models

1. In Bayesian linear regression we place a prior over weights  $\mathbf{w}$  and obtain a posterior  $p(\mathbf{w} \mid \mathbf{y}, \mathbf{X})$ . Explain in words what it means to instead place a prior directly over functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ .
2. A Gaussian process is defined as: a collection of random variables such that every finite subset has a joint Gaussian distribution. Let  $f \sim \mathcal{GP}(m, k)$  where  $m(\mathbf{x}) = 0$  and  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2\right)$ .
  - Write the joint distribution  $p(f(\mathbf{x}_1), f(\mathbf{x}_2))$  for two inputs  $\mathbf{x}_1, \mathbf{x}_2$ .
  - What is the marginal distribution of  $f(\mathbf{x}_1)$ ?

## Kernel Functions

1. The squared exponential (RBF) kernel is

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right).$$

- What does the lengthscale  $\ell$  control? What happens as  $\ell \rightarrow 0$  and  $\ell \rightarrow \infty$ ?
  - What does the signal variance  $\sigma_f^2$  control?
2. For inputs  $x_1 = 0, x_2 = 1, x_3 = 3$  with  $\sigma_f^2 = 1$  and  $\ell = 1$ , compute the  $3 \times 3$  kernel matrix  $\mathbf{K}$  where  $K_{ij} = k(x_i, x_j)$ .
  3. A valid kernel must be symmetric and positive semi-definite. Verify that the RBF kernel is symmetric. Why is positive semi-definiteness important for the covariance matrix of a GP?
  4. Consider the linear kernel  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ . What kind of functions does a GP with this kernel represent?

## GP Regression (Exact Inference)

Consider a noise-corrupted observation model:

$$y_n = f(\mathbf{x}_n) + \varepsilon_n, \quad \varepsilon_n \sim \mathcal{N}(0, \sigma_n^2),$$

with prior  $f \sim \mathcal{GP}(0, k)$ .

1. Write the joint distribution of the observations  $\mathbf{y} = (y_1, \dots, y_N)^\top$  and the latent function values  $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^\top$ .

2. The joint distribution of training outputs  $\mathbf{y}$  and test function values  $\mathbf{f}_\star = f(\mathbf{X}_\star)$  is Gaussian:

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_\star \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \mathbf{K}_\star^\top \\ \mathbf{K}_\star & \mathbf{K}_{\star\star} \end{pmatrix}\right),$$

where  $\mathbf{K} = k(\mathbf{X}, \mathbf{X})$ ,  $\mathbf{K}_\star = k(\mathbf{X}_\star, \mathbf{X})$ ,  $\mathbf{K}_{\star\star} = k(\mathbf{X}_\star, \mathbf{X}_\star)$ . Derive the posterior predictive mean and covariance:

$$\mu_\star = \mathbf{K}_\star(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}, \quad \Sigma_\star = \mathbf{K}_{\star\star} - \mathbf{K}_\star(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{K}_\star^\top.$$

3. Suppose we have a single training point  $x_1 = 0$  with  $y_1 = 1$ , and we want to predict at  $x_\star = 1$ . Use the RBF kernel with  $\sigma_f^2 = 1$ ,  $\ell = 1$ , and  $\sigma_n^2 = 0.1$ .
- Compute  $K = k(x_1, x_1)$ ,  $k_\star = k(x_\star, x_1)$ ,  $k_{\star\star} = k(x_\star, x_\star)$ .
  - Compute the posterior predictive mean  $\mu_\star$  and variance  $\sigma_\star^2$ .
4. **Two training points.** Now suppose we have two training points  $x_1 = 0$ ,  $y_1 = 1$  and  $x_2 = 2$ ,  $y_2 = -1$ , and we want to predict at  $x_\star = 1$ . Use the RBF kernel with  $\sigma_f^2 = 1$ ,  $\ell = 1$ ,  $\sigma_n^2 = 0.25$ .
- Build the  $2 \times 2$  matrix  $\mathbf{K} + \sigma_n^2 \mathbf{I}$  and compute its inverse.
  - Compute the posterior predictive mean  $\mu_\star$  and variance  $\sigma_\star^2$  at  $x_\star = 1$ .
5. **Effect of noise.** Repeat the single-training-point prediction from the previous section ( $x_1 = 0$ ,  $y_1 = 1$ ,  $x_\star = 1$ , RBF with  $\sigma_f^2 = 1$ ,  $\ell = 1$ ) for two noise levels:  $\sigma_n^2 = 0.01$  and  $\sigma_n^2 = 4$ .
- Compute  $\mu_\star$  and  $\sigma_\star^2$  for each case.
  - Describe qualitatively how increasing  $\sigma_n^2$  affects the posterior mean and variance at the test point, and give an intuitive explanation.

## Hyperparameter Optimisation

1. The marginal likelihood (model evidence) for GP regression is

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^\top (\mathbf{K}_\theta + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_\theta + \sigma_n^2 \mathbf{I}| - \frac{N}{2} \log 2\pi.$$

Identify the three terms and give an intuitive explanation of the role each plays in selecting the kernel hyperparameters  $\boldsymbol{\theta} = (\ell, \sigma_f^2, \sigma_n^2)$ .

2. Why is maximising the marginal likelihood preferable to maximising the training data likelihood  $p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$  for choosing hyperparameters?
3. What is the computational complexity of exact GP regression in terms of the number of training points  $N$ ? What bottleneck causes this? Name one strategy to reduce this cost.